

# Static and Dynamic Semantics of NoSQL Languages

Véronique Benzaken<sup>1</sup>   Giuseppe Castagna<sup>2</sup>   Kim Nguyễn<sup>1</sup>   Jérôme Siméon<sup>3</sup>

<sup>1</sup>LRI, Université Paris-Sud, Orsay, France,   <sup>2</sup>CNRS, PPS, Univ Paris Diderot, Sorbonne Paris Cité, Paris, France

<sup>3</sup>IBM Watson Research, Hawthorne, NY, USA

## Abstract

We present a calculus for processing semistructured data that spans differences of application area among several novel query languages, broadly categorized as “NoSQL”. This calculus lets users define their own operators, capturing a wider range of data processing capabilities, whilst providing a typing precision so far typical only of primitive hard-coded operators. The type inference algorithm is based on semantic type checking, resulting in type information that is both precise, and flexible enough to handle structured and semistructured data. We illustrate the use of this calculus by encoding a large fragment of Jaql, including operations and iterators over JSON, embedded SQL expressions, and co-grouping, and show how the encoding directly yields a typing discipline for Jaql as it is, namely without the addition of any type definition or type annotation in the code.

## 1. Introduction

The emergence of Cloud computing, and the ever growing importance of data in applications, has given birth to a whirlwind of new data models [19, 24] and languages. Whether they are developed under the banner of “NoSQL” [12, 35], for BigData Analytics [5, 18, 28], for Cloud computing [3], or as domain specific languages (DSL) embedded in a host language [21, 27, 32], most of them share a common subset of SQL and the ability to handle semistructured data. While there is no consensus yet on the precise boundaries of this class of languages, they all share two common traits: (i) an emphasis on sequence operations (eg, through the popular MapReduce paradigm) and (ii) a lack of types for both data and programs (contrary to, say, XML programming or relational databases where data schemas are pervasive). In [21, 22], Meijer argues that such languages can greatly benefit from formal foundations, and suggests comprehensions [7, 33, 34] as a unifying model. Although we agree with Meijer for the need to provide unified, formal foundations to those new languages, we argue that such foundations should account for novel features critical to various application domains that are not captured by comprehensions. Also, most of those languages provide limited type checking, or ignore it altogether. We believe type checking is essential for many applications, with usage ranging from error detection to optimization. But we understand the designers and programmers of those languages who are averse to any kind of type definition or annotation. In this paper, we propose a calculus which is expressive enough to capture languages that go beyond SQL or comprehensions. We show how the calculus adapts to various data models while retaining a precise type checking that can exploit in a flexible way limited type information, information

that is deduced directly from the structure of the program even in the absence of any explicit type declaration or annotation.

**Example.** We use Jaql [5, 18], a language over JSON [19] developed for BigData analytics, to illustrate how our proposed calculus works. Our reason for using Jaql is that it encompasses all the features found in the previously cited query languages and includes a number of original ones, as well. Like Pig [28] it supports sequence iteration, filtering, and grouping operations on non-nested queries. Like AQL [3] and XQuery [6], it features nested queries. Furthermore, Jaql uses a rich data model that allows arbitrary nesting of data (it works on generic sequences of JSON records whose fields can contain other sequences or records) while other languages are limited to flat data models, such as AQL whose data-model is similar to the standard relational model used by SQL databases (tuples of scalars and of lists of scalars). Lastly, Jaql includes SQL as an embedded sub-language for relational data. For these reasons, although in the present work we focus almost exclusively on Jaql, we believe that our work can be adapted without effort to a wide array of sequence processing languages.

The following Jaql program illustrates some of those features. It performs co-grouping [28] between one JSON input, containing information about departments, and one relational input containing information about employees. The query returns for each department its name and id, from the first input, and the number of high-income employees from the second input. A SQL expression is used to select the employees with income above a given value, while a Jaql filter is used to access the set of departments and the elements of these two collections are processed by the group expression (in Jaql “\$” denotes the current element).

```
group
  (depts -> filter each x (x.size > 50))
    by g = $.deptid as ds,
  (SELECT * FROM employees WHERE income > 100)
    by g = $.dept as es
into { dept: g,
      deptName: ds[0].name,
      numEmps: count(es) };
```

The query blends Jaql expressions (eg, `filter` which selects, in the collection `depts`, departments with a `size` of more than 50 employees, and the grouping itself) with a SQL statement (selecting employees in a relational table for which the salary is more than 100). Relations are naturally rendered in JSON as collections of records. In our example, one of the key difference is that field access in SQL requires the field to be present in the record, while the same operation in Jaql does not. Actually, field selection in Jaql is very expressive since it can be applied also to collections with the effect that the selection is recursively applied to the components of the collection and the collection of the results returned, and similarly for `filter` and other iterators. In other words, the expression

An extended abstract of this work is included in the proceeding of *POPL 13, 40th ACM Symposium on Principles of Programming Languages*, ACM Press, 2013.

`filter each x (x.size > 50)` above will work as much when `x` is bound to a record (with or without a `size` field: in the latter case the selection returns `null`), as when `x` is bound to a collection of records or of arbitrary nested collections thereof. This accounts for the semistructured nature of JSON compared to the relational model. Our calculus can express both, in a way that illustrates the difference in both the dynamic semantics and static typing.

In our calculus, the selection of all records whose *mandatory* field `income` is greater than 100 is defined as:

```
let Sel =
  'nil => 'nil
  | ({income: x, .. } as y, tail) =>
    if x > 100 then (y, Sel(tail)) else Sel(tail)
```

(collections are encoded as lists *à la* Lisp) while the filtering among records or arbitrary nested collections of records of those where the (optional) `size` field is present and larger than 50 is:

```
let Fil =
  'nil => 'nil
  | ({size: x, .. } as y, tail) =>
    if x > 50 then (y, Fil(tail)) else Fil(tail)
  | ((x, xs), tail) => (Fil(x, xs), Fil(tail))
  | (_, tail) => Fil(tail)
```

The terms above show nearly all the basic building blocks of our calculus (only composition is missing), building blocks that we dub *filters*. Filters can be defined recursively (eg, `Sel(tail)` is a recursive call); they can perform pattern matching as found in functional languages (the filter  $p \Rightarrow f$  executes  $f$  in the environment resulting from the matching of pattern  $p$ ); they can be composed in alternation ( $f_1|f_2$  tries to apply  $f_1$  and if it fails it applies  $f_2$ ), they can spread over the structure of their argument (eg,  $(f_1, f_2)$ —of which  $(x, \text{Sel}(\text{tail}))$  is an instance—requires an argument of a product type and applies the corresponding  $f_i$  component-wise).

For instance, the filter `Fil` scans collections encoded as lists *à la* Lisp (ie, by right associative pairs with `'nil` denoting the empty list). If its argument is the empty list, then it returns the empty list; if it is a list whose head is a record with a `size` field (and possibly other fields matched by `..`), then it captures the whole record in  $y$ , the content of the field in  $x$ , the tail of the list in `tail`, and keeps or discards  $y$  (ie, the record) according to whether  $x$  (ie, the field) is larger than 50; if the head is also a list, then it recursively applies both on the head and on the tail; if the head of the list is neither a list, nor a record with a `size` field, then the head is discarded. The encoding of the whole grouping query is given in Section 5.3.

Our aim is not to propose yet another “NoSQL/cloud computing/bigdata analytics” query language, but rather to show how to *express* and *type* such languages via an encoding into our core calculus. Each such language can in this way preserve its execution model but obtain for free a formal semantics, a type inference system and, as it happens, a prototype implementation. The type information is deduced via the encoding (without the need of *any* type annotation) and can be used for early error detection and debugging purposes. The encoding also yields an executable system that can be used for rapid prototyping. Both possibilities are critical in most typical usage scenarios of these languages, where deployment is very expensive both in time and in resources. As observed by Meijer [21] the advent of big data makes it more important than ever for programmers (and, we add, for language and system designers) to have a single abstraction that allows them to process, transform, query, analyze, and compute across data presenting utter variability both in volume and in structure, yielding a “mind-blowing number

For the sake of precision, to comply with Jaql’s semantics the last pattern should rather be  $(\{.. \}, \text{tail}) \Rightarrow \text{Fil}(\text{tail})$ : since field selection  $e.\text{size}$  fails whenever  $e$  is not a record or a list, this definition would detect the possibility of this failure by a static type error.

of new data models, query languages, and execution fabrics” [21]. The framework we present here, we claim, encompasses them all. A long-term goal is that the compilers of these languages could use the type information inferred from the encoding and the encoding itself to devise further optimizations.

**Types.** Pig [28], Jaql [18, 29], AQL [3] have all been conceived by considering just the map-reduce execution model. The type (or, schema) of the manipulated data did not play any role in their design. As a consequence these languages are untyped and, when present, types are optional and clearly added as an afterthought. Differences in data model or type discipline are particularly important when embedded in a host language (since they yield the so-called impedance mismatch). The reason why types were/are disregarded in such languages may originate in an alleged tension between type inference and heterogeneous/semistructured data: on the one hand these languages are conceived to work with collections of data that are weakly or partially structured, on the other hand current languages with type inference (such as Haskell or ML) can work only on homogeneous collections (typically, lists of elements of the same type).

In this work we show that the two visions can coexist: we type data by semantic subtyping [16], a type system conceived for semistructured data, and describe computations by our *filters* which are untyped combinators that, thanks to a technique of weak typing introduced in [9], can polymorphically type the results of data query and processing with a high degree of precision. The conception of *filters* is driven by the schema of the data rather than the execution model and we use them (i) to capture and give a uniform semantics to a wide range of semi structured data processing capabilities, (ii) to give a type system that encompasses the types defined for such languages, if any, notably Pig, Jaql and AQL (but also XML query and processing languages: see Section 5.3), (iii) to infer the precise result types of queries written in these languages *as they are* (so without the addition of any explicit type annotation/definition or new construct), and (iv) to show how minimal extensions/modifications of the current syntax of these languages can bring dramatic improvements in the precision of the inferred types.

The types we propose here are extensible record types and heterogeneous lists whose content is described by regular expressions on types as defined by the following grammar:

Types	$t ::=$	$v$	
		$v$	(singleton)
		$\{\ell:t, \dots, \ell:t\}$	(closed record)
		$\{\ell:t, \dots, \ell:t, ..\}$	(open record)
		$[r]$	(sequences)
		$\text{int} \mid \text{char}$	(base)
		$\text{any} \mid \text{empty} \mid \text{null}$	(special)
		$t t$	(union)
		$t \setminus t$	(difference)

**Regexp**  $r ::= \varepsilon \mid t \mid r* \mid r+ \mid r? \mid rr \mid r|r$

where  $\varepsilon$  denotes the empty word. The semantics of types can be expressed in terms of sets of *values* (values are either constants—such as 1, 2, `true`, `false`, `null`, `'1'`, the latter denoting the character 1—, records of values, or lists of values). So the singleton type  $v$  is the type that contains just the value  $v$  (in particular `null` is the singleton type containing the value `null`). The closed record type  $\{a:\text{int}, b:\text{int}\}$  contains all record values with exactly two fields  $a$  and  $b$  with integer values, while the open record type  $\{a:\text{int}, b:\text{int}, ..\}$  contains all record values with at least two fields  $a$  and  $b$  with integer values. The sequence type  $[r]$  is the set of all sequences whose content is described by the *regular expression*  $r$ ; so, for example  $[\text{char}^*]$  contains all sequences of characters (we will use `string` to denote this type and the standard double quote notation to denote its values) while  $[(\{a:\text{int}\} \{a:\text{int}\})^+]$

denotes nonempty lists of even length containing record values of type  $\{a:int\}$ . The union type  $s|t$  contains all the values of  $s$  and of  $t$ , while the difference type  $s\backslash t$  contains all the values of  $s$  that are not in  $t$ . We shall use `bool` as an abbreviation of the union of the two singleton types containing true and false: `'true'|'false'`. `any` and `empty` respectively contain all and no values. Recursive type definitions are also used (see Section 2.2 for formal details).

These types can express all the types of Pig, Jaql and AQL, all XML types, and much more. So for instance, AQL includes only homogeneous lists of type  $t$ , that can be expressed by our types as  $[t^*]$ . In Jaql's documentation one can find the type  $[long(value=1), string(value="a"), boolean^*]$  which is the type of arrays whose first element is 1, the second is the string "a" and all the other are booleans. This can be easily expressed in our types as  $[1\ "a"\ bool^*]$ . But while Jaql only allows a limited use of regular expressions (Kleene star can only appear in tail position) our types do not have such restrictions. So for example  $[char^* '@' char^* '.' ((f' r')|(i' t'))]$  is the type of all strings (*ie*, sequences of chars) that denote email addresses ending by either `.fr` or `.it`. We use some syntactic sugar to make terms as the previous one more readable (*eg*,  $[.* '@' .* ('.fr'|'.it')]$ ). Likewise, henceforth we use  $\{a?:t\}$  to denote that the field  $a$  of type  $t$  is optional; this is just syntactic sugar for stating that either the field is undefined or it contains a value of type  $t$  (for the formal definition see Appendix G).

Coming back to our initial example, the filter `Fi1` defined before expects as argument a collection of the following type:

```
type Depts = [ ( {size?: int, ..} | Depts ) * ]
```

that is, a possibly empty, arbitrary nested list of records with an optional `size` field of type `int`: notice that it is important to specify the optional field and its type since a `size` field of a different type would make the expression `x > 50` raise a run-time error. This information is deduced just from the *structure* of the filter (since `Fi1` does not contain any type definition or annotation).

We define a type inference system that rejects any argument of `Fi1` that has not type `Depts`, and deduces for arguments of type  $[(\{size: int, addr: string\} | \{sec: int\} | Depts)^+]$  (which is a subtype of `Depts`) the result type  $[(\{size: int, addr: string\} | Depts)^*]$  (so it does not forget the field `addr` but discards the field `sec`, and by replacing `*` for `+` recognizes that the test may fail).

By encoding primitive Jaql operations into a formal core calculus we shall provide them a formal and clean semantics as well as precise typing. So for instance it will be clear that applying the following dot selection  $[(\{a:3\} \{a:5, b:true\})].a$  the result will be  $[ [3] \ 5 ]$  and we shall be able to deduce that `_.a` applied to arbitrary nested lists of records with an optional integer  $a$  field (*ie*, of type  $t = \{a?:int\} | [t^*]$ ) yields arbitrary nested lists of `int` or `null` values (*ie*, of type  $u = int | null | [u^*]$ ).

Finally we shall show that if we accept to extend the current syntax of Jaql (or of some other language) by some minimal filter syntax (*eg*, the pattern filter) we can obtain a huge improvement in the precision of type inference.

**Contributions.** The main contribution of this work is the definition of a calculus that encompasses structural operators scattered over NoSQL languages and that possesses some characteristics that make it unique in the swarm of current semi-structured data processing languages. In particular it is parametric (though fully embeddable) in a host language; it uniformly handles both width and deep nested data recursion (while most languages offer just the

former and limited forms of the latter); finally, it includes first-class arbitrary deep composition (while most languages offer this operator only at top level), whose power is nevertheless restrained by the type system.

An important contribution of this work is that it directly compares a programming language approach with the tree transducer one. Our calculus implements transformations typical of top-down tree transducers but has several advantages over the transducer approach: (1) the transformations are expressed in a formalism immediately intelligible to any functional programmer; (2) our calculus, in its untyped version, is Turing complete; (3) its transformations can be statically typed (at the expenses of Turing completeness) without any annotation yielding precise result types (4) even if we restrict the calculus only to well-typed terms (thus losing Turing completeness), it still is strictly more expressive than well-known and widely studied deterministic top-down tree transducer formalisms.

The technical contributions are (i) the proof of Turing completeness for our formalism, (ii) the definition of a type system that copes with records with computable labels (iii) the definition of a static type system for filters and its correctness, (iv) the definition of a static analysis that ensures the termination (and the proof thereof) of the type inference algorithm with complexity bounds expressed in the size of types and filters and (v) the proof that the terms that pass the static analysis form a language strictly more expressive than top-down tree transducers.

**Outline.** In Section 2 we present the syntax of the three components of our system. Namely, a minimal set of *expressions*, the calculus of *filters* used to program user-defined operators or to encode the operators of other languages, and the core *types* in which the types we just presented are to be encoded. Section 3 defines the operational semantics of filters and a declarative semantics for operators. The type system as well as the type inference algorithm are described in Section 4. In Section 5 we present how to handle a large subset of Jaql. Section 8 reports on some subtler design choices of our system. We compare related works in Section 9 and conclude in Section 10. In order to avoid blurring the presentation, proofs, secondary results, further encodings, and extensions are moved into a separate appendix.

## 2. Syntax

In this section we present the syntax of the three components of our system: a minimal set of *expressions*, the calculus of *filters* used to program user-defined operators or to encode the operators of other languages, and the core *types* in which the types presented in the introduction are to be encoded.

The core of our work is the definition of filters and types. The key property of our development is that filters can be grafted to any host language that satisfies minimal requirements, by simply adding filter application to the expressions of the host language. The minimal requirements of the host language for this to be possible are quite simple: it must have constants (typically for types `int`, `char`, `string`, and `bool`), variables, and either pairs or record values (not necessarily both). On the static side the host language must have at least basic and products types and be able to assign a type to expressions in a given type environment (*ie*, under some typing assumptions for variables). By the addition of filter applications, the host language can acquire or increase the capability to define polymorphic user-defined iterators, query and processing expressions, and be enriched with a powerful and precise type system.

### 2.1 Expressions

In this work we consider the following set of expressions

The only exception are the “bags” types we did not include in order to focus on essential features.



### Definition 1 (expressions).

<b>Exprs</b>	$e ::= c$	(constants)
	$x$	(variables)
	$(e, e)$	(pairs)
	$\{e:e, \dots, e:e\}$	(records)
	$e + e$	(record concatenation)
	$e \setminus \ell$	(field deletion)
	$\text{op}(e, \dots, e)$	(built-in operators)
	$f e$	(filter application)

where  $f$  ranges over *filters* (defined later on),  $c$  over generic constants, and  $\ell$  over *string* constants.

Intuitively, these expressions represent the syntax supplied by the host language —though only the first two and one of the next two are really needed— that we extend with (the missing expressions and) the expression of filter application. Expressions are formed by constants, variables, pairs, records, and operation on records: record concatenation gives priority to the expression on the right. So if in  $r_1 + r_2$  both records contains a field with the same label, it is the one in  $r_2$  that will be taken, while field deletion does not require the record to contain a field with the given label (though this point is not important). The metavariable  $\text{op}$  ranges over operators as well as functions and other constructions belonging to or defined by the host language. Among expressions we single out a set of *values*, intuitively the results of computations, that are formally defined as follows:

$$v ::= c \mid (v, v) \mid \{\ell:v; \dots; \ell:v\}$$

We use "foo" for character string constants, 'c' for characters, 1 2 3 4 5 and so on for integers, and backquoted words, such as 'foo, for atoms (*ie*, user-defined constants). We use three distinguished atoms 'nil', 'true, and 'false. Double quotes can be omitted for strings that are labels of record fields: thus we write {name:"John"} rather than {"name":"John"}. Sequences (*aka*, heterogeneous lists, ordered collections, arrays) are encoded *à la LISP*, as nested pairs where the atom 'nil denotes the empty list. We use  $[e_1 \dots e_n]$  as syntactic sugar for  $(e_1, \dots, (e_n, \text{'nil}))$ .

## 2.2 Types

Expressions, in particular filter applications, are typed by the following set of types (typically only basic, product, recursive and —some form of— record types will be provided by the host language):

### Definition 2 (types).

<b>Types</b>	$t ::= b$	(basic types)
	$v$	(singleton types)
	$(t, t)$	(products)
	$\{\ell:t, \dots, \ell:t\}$	(closed records)
	$\{\ell:t, \dots, \ell:t, ..\}$	(open records)
	$t t$	(union types)
	$t \& t$	(intersection types)
	$\neg t$	(negation type)
	$\text{empty}$	(empty type)
	$\text{any}$	(any type)
	$\mu T. t$	(recursive types)
	$T$	(recursion variable)
	$\text{Op}(t, \dots, t)$	(foreign type calls)

where every recursion is guarded, that is, every type variable is separated from its binder by at least one application of a type constructor (*ie*, products, records, or  $\text{Op}$ ).

Most of these types were already explained in the introduction. We have basic types ( $\text{int}$ ,  $\text{bool}$ , ...) ranged over by  $b$  and singleton types  $v$  denoting the type that contains only the value  $v$ .

Record types come in two flavors: closed record types whose values are records with exactly the fields specified by the type, and open record types whose values are records with *at least* the fields specified by the type. Product types are standard and we have a complete set of type connectives, that is, finite unions, intersections and negations. We use  $\text{empty}$ , to denote the type that has no values and  $\text{any}$  for the type of all values (sometimes denoted by " $\_$ " when used in patterns). We added a term for recursive types, which allows us to encode both the regular expression types defined in the introduction and, more generally, the recursive type definitions we used there. Finally, we use  $\text{Op}$  (capitalized to distinguish it from expression operators) to denote the host language's *type* operators (if any). Thus, when filter applications return values whose type belongs just to the foreign language (*eg*, a list of functions) we suppose the typing of these functions be given by some type operators. For instance, if  $\text{succ}$  is a user defined successor function, we will suppose to be given its type in the form  $\text{Arrow}(\text{int}, \text{int})$  and, similarly, for its application, say  $\text{apply}(\text{succ}, 3)$  we will be given the type of this expression (presumably  $\text{int}$ ). Here  $\text{Arrow}$  is a type operator and  $\text{apply}$  an expression operator.

The denotational semantics of types as sets of values, that we informally described in the introduction, is at the basis of the definition of the subtyping relation for these types. We say that a type  $t_1$  is a subtype of a type  $t_2$ , noted  $t_1 \leq t_2$ , if and only if the set of values denoted by  $t_1$  is contained (in the set-theoretic sense) in the set of values denoted by  $t_2$ . For the formal definition and the decision procedure of this subtyping relation the reader can refer to the work on semantic subtyping [16].

## 2.3 Patterns

Filters are our core untyped operators. All they can do are three different things: (1) they can structurally decompose and transform the values they are applied to, or (2) they can be sequentially composed, or (3) they can do pattern matching. In order to define filters, thus, we first need to define patterns.

### Definition 3 (patterns).

<b>Patterns</b>	$p ::= t$	(type)
	$x$	(variable)
	$(p, p)$	(pair)
	$\{\ell:p, \dots, \ell:p\}$	(closed rec)
	$\{\ell:p, \dots, \ell:p, ..\}$	(open rec)
	$p p$	(or/union)
	$p \& p$	(and/intersection)

where the subpatterns forming pairs, records, and intersections have distinct capture variables, and those forming unions have the same capture variables.

Patterns are essentially types in which capture variables (ranged over by  $x, y, \dots$ ) may occur in every position that is not under a negation or a recursion. A pattern is used to match a value. The matching of a value  $v$  against a pattern  $p$ , noted  $v/p$ , either fails (noted  $\Omega$ ) or it returns a substitution from the variables occurring in the pattern, into values. The substitution is then used as an environment in which some expression is evaluated. If the pattern is a type, then the matching fails if and only if the pattern is matched against a value that does not have that type, otherwise it returns the empty substitution. If it is a variable, then the matching always succeeds and returns the substitution that assigns the matched value to the variable. The pair pattern  $(p_1, p_2)$  succeeds if and only if it is matched against a pair of values and each sub-pattern succeeds on the corresponding projection of the value (the union of the two substitutions is then returned). Both record patterns are similar to the product pattern with the specificity that in the open record pattern " $..$ " matches all the fields that are not specified in the

pattern. An intersection pattern  $p_1 \& p_2$  succeeds if and only if both patterns succeed (the union of the two substitutions is then returned). The union pattern  $p_1 | p_2$  first tries to match the pattern  $p_1$  and if it fails it tries the pattern  $p_2$ .

For instance, the pattern  $(\text{int} \& x, y)$  succeeds only if the matched value is a pair of values  $(v_1, v_2)$  in which  $v_1$  is an integer—in which case it returns the substitution  $\{x/v_1, y/v_2\}$ —and fails otherwise. Finally notice that the notation “ $p$  as  $x$ ” we used in the examples of the introduction, is syntactic sugar for  $p \& x$ .

This informal semantics of matching (see [16] for the formal definition) explains the reasons of the restrictions on capture variables in Definition 3: in intersections, pairs, and records all patterns must be matched and, thus, they have to assign distinct variables, while in union patterns just one pattern will be matched, hence the same set of variables must be assigned, whichever alternative is selected.

The strength of patterns is their connections with types and the fact that the pattern matching operator can be typed *exactly*. This is entailed by the following theorems (both proved in [16]):

**Theorem 4 (Accepted type [16]).** *For every pattern  $p$ , the set of all values  $v$  such that  $v/p \neq \Omega$  is a type. We call this set the accepted type of  $p$  and note it by  $\llbracket p \rrbracket$ .*

The fact that the exact set of values for which a matching succeeds is a type is not obvious. It states that for every pattern  $p$  there exists a syntactic type produced by the grammar in Definition 2 whose semantics is exactly the set of all and only values that are matched by  $p$ . The existence of this syntactic type, which we note  $\llbracket p \rrbracket$ , is of utmost importance for a precise typing of pattern matching. In particular, given a pattern  $p$  and a type  $t$  contained in (ie, subtype of)  $\llbracket p \rrbracket$ , it allows us to compute the *exact* type of the capture variables of  $p$  when it is matched against a value in  $t$ :

**Theorem 5 (Type environment [16]).** *There exists an algorithm that for every pattern  $p$ , and  $t \leq \llbracket p \rrbracket$  returns a type environment  $t/p \in \text{Vars}(p) \rightarrow \text{Types}$  such that  $(t/p)(x) = \{(v/p)(x) \mid v : t\}$ .*

## 2.4 Filters

**Definition 6 (filters).** A *filter* is a term generated by:

<b>Filters</b>	$f$	$::=$	$e$	(expression)
			$p \Rightarrow f$	(pattern)
			$(f, f)$	(product)
			$\{\ell : f, \dots, \ell : f, ..\}$	(record)
			$f   f$	(union)
			$\mu X. f$	(recursion)
			$X a$	(recursive call)
			$f ; f$	(composition)
			$o$	(declarative operators)
<b>Operators</b>	$o$	$::=$	$\text{groupby } f$	(filter grouping)
			$\text{orderby } f$	(filter ordering)
<b>Arguments</b>	$a$	$::=$	$x$	(variables)
			$c$	(constants)
			$(a, a)$	(pairs)
			$\{\ell : a, \dots, \ell : a\}$	(record)

such that for every subterm of the form  $f ; g$ , no recursion variable is free in  $f$ .

Filters are like transducers, that when applied to a value return another value. However, unlike transducers they possess more “programming-oriented” constructs, like the ability to test an input and capture subterms, recompose an intermediary result from captured values and a composition operator. We first describe informally the semantics of each construct.

The expression filter  $e$  always returns the value corresponding to the evaluation of  $e$  (and discards its argument). The filter  $p \Rightarrow f$  applies the filter  $f$  to its argument in the environment obtained by matching the argument against  $p$  (provided that the matching does not fail). This rather powerful feature allows a filter to perform two critical actions: (i) inspect an input with regular pattern-matching before exploring it and (ii) capture part of the input that can be reused during the evaluation of the subfilter  $f$ . If the argument application of  $f_i$  to  $v_i$  returns  $v'_i$  then the application of the product filter  $(f_1, f_2)$  to an argument  $(v_1, v_2)$  returns  $(v'_1, v'_2)$ ; otherwise, if any application fails or if the argument is not a pair, it fails. The record filter is similar: it applies to each specified field the corresponding filter and, as stressed by the “ $..$ ”, leaves the other fields unchanged; it fails if any of the applications does, or if any of the specified fields is absent, or if the argument is not a record. The filter  $f_1 | f_2$  returns the application of  $f_1$  to its argument or, if this fails, the application of  $f_2$ . The semantics of a recursive filter is given by standard unfolding of its definition in recursive calls. The only real restriction that we introduce for filters is that recursive calls can be done only on arguments of a given form (ie, on arguments that have the form of values where variables may occur). This restriction in practice amounts to forbid recursive calls on the result of another recursively defined filter (all other cases can be easily encoded). The reason of this restriction is technical, since it greatly simplifies the analysis of Section 4.5 (which ensures the termination of type inference) without hampering expressiveness: filters are Turing complete even with this restriction (see Theorem 7). Filters can be composed: the filter  $f_1 ; f_2$  applies  $f_2$  to the result of applying  $f_1$  to the argument and fails if any of the two does. The condition that in every subterm of the form  $f ; g$ ,  $f$  does not contain free recursion variables is not strictly necessary. Indeed, we could allow such terms. The point is that the analysis for the termination of the typing would then reject all such terms (apart from trivial ones in which the result of the recursive call is not used in the composition). But since this restriction *does not* restrict the expressiveness of the calculus (Theorem 7 proves Turing completeness with this restriction), then the addition of this restriction is just a design (rather than a technical) choice: we prefer to forbid the programmer to write recursive calls on the left-hand side of a composition, than systematically reject all the programs that use them in a non-trivial way.

Finally, we singled out some specific filters (specifically, we chose `groupby` and `orderby`) whose semantics is generally specified in a declarative rather than operational way. These do not bring any expressive power to the calculus (the proof of Turing completeness, Theorem 7, does not use these declarative operators) and actually they can be encoded by the remaining filters, but it is interesting to single them out because they yield either simpler encodings or more precise typing.

## 3. Semantics

The operational semantics of our calculus is given by the reduction semantics for filter application and for the record operations. Since the former is the only novelty of our work, we save space and omit the latter, which are standard anyhow.

### 3.1 Big step semantics

We define a big step operational semantics for filters. The definition is given by the inference rules in Figure 1 for judgements of the form  $\delta ; \gamma \vdash_{\text{eval}} f(a) \rightsquigarrow r$  and describes how the evaluation of the application of filter  $f$  to an argument  $a$  in an environment  $\gamma$  yields an object  $r$  where  $r$  is either a value or  $\Omega$ . The latter is a special value which represents a runtime error: it is raised by the rule (**error**) either because a filter did not match the form of its argument (eg, the argument of a filter product was not a pair)

<b>(expr)</b>	$\frac{}{\delta; \gamma \vdash_{eval} e(v) \rightsquigarrow r} \quad r = eval(\gamma, e)$	
<b>(prod)</b>	$\frac{\delta; \gamma \vdash_{eval} f_1(v_1) \rightsquigarrow r_1 \quad \delta; \gamma \vdash_{eval} f_2(v_2) \rightsquigarrow r_2}{\delta; \gamma \vdash_{eval} (f_1, f_2)(v_1, v_2) \rightsquigarrow (r_1, r_2)} \quad \begin{array}{l} \text{if } r_1 \neq \Omega \\ \text{and } r_2 \neq \Omega \end{array}$	<b>(union1)</b> $\frac{\delta; \gamma \vdash_{eval} f_1(v) \rightsquigarrow r_1}{\delta; \gamma \vdash_{eval} (f_1   f_2)(v) \rightsquigarrow r_1} \quad \text{if } r_1 \neq \Omega$
<b>(patt)</b>	$\frac{\delta; \gamma, v/p \vdash_{eval} f(v) \rightsquigarrow r}{\delta; \gamma \vdash_{eval} (p \Rightarrow f)(v) \rightsquigarrow r} \quad \text{if } v/p \neq \Omega$	<b>(union2)</b> $\frac{\delta; \gamma \vdash_{eval} f_1(v) \rightsquigarrow \Omega \quad \delta; \gamma \vdash_{eval} f_2(v) \rightsquigarrow r_2}{\delta; \gamma \vdash_{eval} (f_1   f_2)(v) \rightsquigarrow r_2}$
<b>(comp)</b>	$\frac{\delta; \gamma \vdash_{eval} f_1(v) \rightsquigarrow r_1 \quad \delta; \gamma \vdash_{eval} f_2(r_1) \rightsquigarrow r_2}{\delta; \gamma \vdash_{eval} (f_1; f_2)(v) \rightsquigarrow r_2} \quad \text{if } r_1 \neq \Omega$	<b>(rec)</b> $\frac{\delta, (X \mapsto f); \gamma \vdash_{eval} f(v) \rightsquigarrow r}{\delta; \gamma \vdash_{eval} (\mu X. f)(v) \rightsquigarrow r}$
<b>(recd)</b>	$\frac{\delta; \gamma \vdash_{eval} f_1(v_1) \rightsquigarrow r_1 \quad \dots \quad \delta; \gamma \vdash_{eval} f_n(v_n) \rightsquigarrow r_n}{\delta; \gamma \vdash_{eval} \{\ell_1: f_1, \dots, \ell_n: f_n, \dots\}(\{\ell_1: v_1, \dots, \ell_n: v_n, \dots, \ell_{n+k}: v_{n+k}\}) \rightsquigarrow \{\ell_1: r_1, \dots, \ell_n: r_n, \dots, \ell_{n+k}: v_{n+k}\}} \quad \text{if } \forall i, r_i \neq \Omega$	<b>(rec-call)</b> $\frac{\delta; \gamma \vdash_{eval} (\delta(X))(a) \rightsquigarrow r}{\delta; \gamma \vdash_{eval} (Xa)(v) \rightsquigarrow r}$
		<b>(error)</b> $\frac{}{\delta; \gamma \vdash_{eval} f(a) \rightsquigarrow \Omega} \quad \text{if no other rule applies}$

Figure 1. Dynamic semantics of filters

or because some pattern matching failed (*ie*, the side condition of **(patt)** did not hold). Notice that the argument  $a$  of a filter is always a value  $v$  unless the filter is the unfolding of a recursive call, in which case variables may occur in it (*cf.* rule **rec-call**). Environment  $\delta$  is used to store the body of recursive definitions.

The semantics of filters is quite straightforward and inspired by the semantics of patterns. The *expression* filter discards its input and evaluates (rather, asks the host language to evaluate) the expression  $e$  in the current environment (**expr**). It can be thought of as the right-hand side of a branch in a `match_with` construct.

The *product* filter expects a pair as input, applies its sub-filters component-wise and returns the pair of the results (**prod**). This filter is used in particular to express sequence mapping, as the first component  $f_1$  transforms the element of the list and  $f_2$  is applied to the tail. In practice it is often the case that  $f_2$  is a recursive call that iterates on arbitrary lists and stops when the input is `nil`. If the input is not a pair, then the filter fails (rule **error**) applies).

The *record* filter expects as input a record value with *at least* the same fields as those specified by the filter. It applies each sub-filter to the value in the corresponding field leaving the contents of other fields unchanged (**recd**). If the argument is not a record value or it does not contain all the fields specified by the record filter, or if the application of any subfilter fails, then the whole application of the record filter fails.

The *pattern* filter matches its input value  $v$  against the pattern  $p$ . If the matching fails so the filter does, otherwise it evaluates its sub-filter in the environment augmented by the substitution  $v/p$  (**patt**).

The *alternative* filter follows a standard first-match policy: If the filter  $f_1$  succeeds, then its result is returned (**union-1**). If  $f_1$  fails, then  $f_2$  is evaluated against the input value (**union-2**). This filter is particularly useful to write the alternative of two (or more) *pattern* filters, making it possible to conditionally continue a computation based on the shape of the input.

The *composition* allows us to pass the result of  $f_1$  as input to  $f_2$ . The composition filter is of paramount importance. Indeed, without it, our only way to iterate (deconstruct) an input value is to use a *product* filter, which always rebuilds a pair as result.

Finally, a *recursive* filter is evaluated by recording its body in  $\delta$  and evaluating it (**rec**), while for a *recursive call* we replace the recursion variable by its definition (**rec-call**).

This concludes the presentation of the semantics of non-declarative filters (*ie*, without groupby and orderby). These form a Turing complete formalism (full proof in Appendix B):

**Theorem 7 (Turing completeness).** *The language formed by constants, variables, pairs, equality, and applications of non-declarative filters is Turing complete.*

*Proof (sketch).* We can encode untyped call-by-value  $\lambda$ -calculus by first applying continuation passing style (CPS) transformations and encoding CPS term reduction rules and substitutions via filters. Thanks to CPS we eschew the restrictions on composition.  $\square$

### 3.2 Semantics of declarative filters

To conclude the presentation of the semantics we have to define the semantics of groupby and orderby. We prefer to give the semantics in a declarative form rather than operationally in order not to tie it to a particular order (of keys or of the execution):

**Groupby:** groupby  $f$  applied to a sequence  $[v_1 \dots v_m]$  reduces to a sequence  $[(k_1, l_1) \dots (k_n, l_n)]$  such that:

1.  $\forall i, 1 \leq i \leq m, \exists j, 1 \leq j \leq n, \text{ s.t. } k_j = f(v_i)$
2.  $\forall j, 1 \leq j \leq n, \exists i, 1 \leq i \leq m, \text{ s.t. } k_j = f(v_i)$
3.  $\forall j, 1 \leq j \leq n, l_j$  is a sequence:  $[v_1^j \dots v_{n_j}^j]$
4.  $\forall j, 1 \leq j \leq n, \forall k, 1 \leq k \leq n_j, f(v_k^j) = k_j$
5.  $k_i = k_j \Rightarrow i = j$
6.  $l_1, \dots, l_n$  is a partition of  $[v_1 \dots v_m]$

**Orderby:** orderby  $f$  applied to  $[v_1 \dots v_n]$  reduces to  $[v'_1 \dots v'_n]$  such that:

1.  $[v'_1 \dots v'_n]$  is a permutation of  $[v_1 \dots v_n]$ ,
2.  $\forall i, j \text{ s.t. } 1 \leq i \leq j \leq n, f(v_i) \leq f(v_j)$

Since the semantics of both operators is deeply connected to a notion of equality and order on values of the host language, we give them as “built-in” operations. However we will illustrate how our type algebra allows us to provide very precise typing rules, specialized for their particular semantics. It is also possible to encode co-grouping (or groupby on several input sequences) with a combination of groupby and filters (*cf.* Appendix H).

### 3.3 Syntactic sugar

Now that we have formally defined the semantics of filters we can use them to introduce some syntactic sugar.

**Expressions.** The reader may have noticed that the productions for expressions (Definition 1) do not define any destructor (*eg*, projections, label selection, ...), just constructors. The reason is that destructors, as well as other common expressions, can be encoded by filter applications:

$$\begin{aligned}
 e.l &\stackrel{\text{def}}{=} (\{\ell: x, \dots\} \Rightarrow x)e \\
 \text{fst}(e) &\stackrel{\text{def}}{=} ((x, \text{any}) \Rightarrow x)e \\
 \text{snd}(e) &\stackrel{\text{def}}{=} ((\text{any}, x) \Rightarrow x)e \\
 \text{let } p = e_1 \text{ in } e_2 &\stackrel{\text{def}}{=} (p \Rightarrow e_2)e_1 \\
 \text{if } e \text{ then } e_1 \text{ else } e_2 &\stackrel{\text{def}}{=} ('true \Rightarrow e_1 | 'false \Rightarrow e_2)e
 \end{aligned}$$

$\text{match } e \text{ with } p_1 \Rightarrow e_1 \mid \dots \mid p_n \Rightarrow e_n$   
 $\stackrel{\text{def}}{=} (p_1 \Rightarrow e_1 \mid \dots \mid p_n \Rightarrow e_n)e$

These are just a possible choice, but others are possible. For instance in Jaql dot selection is overloaded: when  $_.\ell$  is applied to a record, Jaql returns the content of its  $\ell$  field; if the field is absent or the argument is `null`, then Jaql returns `null` and fails if the argument is not a record; when applied to a list (‘array’ in Jaql terminology) it recursively applies to all the elements of the list. So Jaql’s “ $_.\ell$ ” is precisely defined as

$\mu X.(\{\ell:x, \dots\} \Rightarrow x \mid (\{\dots\} \mid \text{null}) \Rightarrow \text{null} \mid (h,t) \Rightarrow (Xh, Xt))$

Besides the syntactic sugar above, in the next section we will use  $t_1 + t_2$  to denote the record type formed by all field types in  $t_2$  and all the field types in  $t_1$  whose label is not already present in  $t_2$ . Similarly  $t \setminus \ell$  will denote the record types formed by all field types in  $t$  apart from the one labelled by  $\ell$ , if present. Finally, we will also use for expressions, types, and patterns the syntactic sugar for lists used in the introduction. So, for instance,  $[p_1 \ p_2 \ \dots \ p_n]$  is matched by lists of  $n$  elements provided that their  $i$ -th element matches  $p_i$ .

## 4. Type inference

In this section we describe a type inference algorithm for our expressions.

### 4.1 Typing of simple and foreign expressions

Variables, constants, and pairs are straightforwardly typed by

[VARS]	[CONSTANT]	[PROD]
$\frac{}{\Gamma \vdash x : \Gamma(x)}$	$\frac{}{\Gamma \vdash c : c}$	$\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash (e_1, e_2) : (t_1, t_2)}$

where  $\Gamma$  denotes a typing environment that is a function from expression variables to types and  $c$  denotes both a constant and the singleton type containing the constant. Expressions of the host language are typed by the `type` function which given a type environment and a foreign expression returns the type of the expression, and that we suppose to be given for each host language.

[FOREIGN]
$\frac{\Gamma \vdash e_1 : t_1 \quad \dots \quad \Gamma \vdash e_n : t_n}{\Gamma \vdash \text{op}(e_1, \dots, e_n) : \text{type}((\Gamma, x_1:t_1, \dots, x_n:t_n), \text{op}(x_1, \dots, x_n))}$

Since the various  $e_i$  can contain filter applications, thus unknown to the host language’s type system, the rule [FOREIGN] swaps them with variables having the same type.

Notice that our expressions, whereas they include *filter* applications, do not include applications of expressions to expressions. Therefore if the host language provides function definitions, then the applications of the host language must be dealt as foreign expressions, as well (*cf.* the expression operator apply in Section 2.2).

### 4.2 Typing of records

The typing of records is novel and challenging because record expressions may contain string *expressions* in label position, such as in  $\{e_1:e_2\}$ , while in all type systems for record we are aware of, labels are never computed. It is difficult to give a type to  $\{e_1:e_2\}$  since, in general, we do not statically know the value that  $e_1$  will return, and which is required to form a record type. All we can (and must) ask is that this value will be a string. To type a record expression  $\{e_1:e_2\}$ , thus, we distinguish two cases according to

The function type must be able to handle type environments with types of our system. It can do it either by subsuming variable with specific types to the types of the host language (*eg.* if the host language does not support singleton types then the singleton type 3 will be subsumed to `int`) or by typing foreign expressions by using our types.

whether the type  $t_1$  of  $e_1$  is finite (*ie.* it contains only finitely many values, such as, say, `Bool`) or not. If a type is finite, (finiteness of regular types seen as tree automata can be decided in polynomial time [10]), then it is possible to write it as a finite union of values (actually, of singleton types). So consider again  $\{e_1:e_2\}$  and let  $t_1$  be the type of  $e_1$  and  $t_2$  the type of  $e_2$ . First,  $t_1$  must be a subtype of `string` (since record labels are strings). So if  $t_1$  is finite it can be expressed as  $\ell_1 \mid \dots \mid \ell_n$  which means that  $e_1$  will return the string  $\ell_i$  for some  $i \in [1..n]$ . Therefore  $\{e_1:e_2\}$  will have type  $\{\ell_i : t_2\}$  for some  $i \in [1..n]$  and, thus, the union of all these types, as expressed by the rule [RCD-FIN] below. If  $t_1$  is infinite instead, then all we can say is that it will be a record with some (unknown) labels, as expressed by rule [RCD-INF].

$$\frac{[\text{RCD-FIN}] \quad \Gamma \vdash e : \ell_1 \mid \dots \mid \ell_n \quad \Gamma \vdash e' : t}{\Gamma \vdash \{e:e'\} : \{\ell_1:t\} \mid \dots \mid \{\ell_n:t\}}$$

$$\frac{[\text{RCD-INF}] \quad \Gamma \vdash e : t \quad \Gamma \vdash e' : t' \quad t \leq \text{string} \quad t \text{ is infinite}}{\Gamma \vdash \{e:e'\} : \{..\}}$$

$$\frac{[\text{RCD-MUL}] \quad \Gamma \vdash \{e_1:e'_1\} : t_1 \quad \dots \quad \Gamma \vdash \{e_n:e'_n\} : t_n}{\Gamma \vdash \{e_1:e'_1, \dots, e_n:e'_n\} : t_1 + \dots + t_n}$$

$$\frac{[\text{RCD-CONC}] \quad \Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2 \quad t_1 \leq \{..\} \quad t_2 \leq \{..\}}{\Gamma \vdash e_1 + e_2 : t_1 + t_2} \quad \frac{[\text{RCD-DEL}] \quad \Gamma \vdash e : t}{\Gamma \vdash e \setminus \ell : t \setminus \ell} \quad t \leq \{..\}$$

Records with multiple fields are handled by the rule [RCD-MUL] which “merges” the result of typing single fields by using the type operator  $+$  as defined in CDuce [4, 15], which is a right-priority record concatenation defined to take into account undefined and unknown fields: for instance,  $\{a:\text{int}, b:\text{int}\} + \{a?:\text{bool}\} = \{a:\text{int} \mid \text{bool}, b:\text{int}\}$ ; unknown fields in the right-hand side may override known fields of the left-hand side, which is why, for instance, we have  $\{a:\text{int}, b:\text{bool}\} + \{b:\text{int}, \dots\} = \{b:\text{int}, \dots\}$ ; likewise, for every record type  $t$  (*ie.* for every  $t$  subtype of  $\{..\}$ ) we have  $t + \{..\} = \{..\}$ . Finally, [RCD-CONC] and [RCD-DEL] deal with record concatenation and field deletion, respectively, in a straightforward way: the only constraint is that all expressions must have a record type (*ie.* the constraints of the form  $\dots \leq \{..\}$ ). See Appendix G for formal definitions of all these type operators.

Notice that these rules do not ensure that a record will not have two fields with the same label, which is a run-time error. Detecting such an error needs sophisticated type systems (*eg.* dependent types) beyond the scope of this work. This is why in the rule [RCD-MUL] we used type operator “ $+$ ” which, in case of multiple occurring labels, since records are unordered, corresponds to randomly choosing one of the types bound to these labels: if such a field is selected, it would yield a run-time error, so its typing can be ambiguous. We can fine tune the rule [RCD-MUL] so that when all the  $t_i$  are finite unions of record types, then we require to have pairwise disjoint sets of labels; but since the problem would still persist for infinite types we prefer to retain the current, simpler formulation.

### 4.3 Typing of filter application

Filters are not first-class: they can be applied but not passed around or computed. Therefore we do not assign types to filters but, as for any other expression, we assign types to *filter applications*. The typing rule for filter application

$$\frac{[\text{FILTER-APP}] \quad \Gamma \vdash e : t \quad \Gamma ; \emptyset ; \emptyset \vdash_{\text{fil}} f(t) : s}{\Gamma \vdash fe : s}$$



relies on an auxiliary deduction system for judgments of the form  $\Gamma; \Delta; M \vdash_{\mu} f(t) : s$  that states that if in the environments  $\Gamma, \Delta, M$  (explained later on) we apply the filter  $f$  to a value of type  $t$ , then it will return a result of type  $s$ .

To define this auxiliary deduction system, which is the core of our type analysis, we first need to define  $\lfloor f \rfloor$ , the type accepted by a filter  $f$ . Intuitively, this type gives a necessary condition on the input for the filter not to fail:

**Definition 8 (Accepted type).** Given a filter  $f$ , the *accepted type* of  $f$ , written  $\lfloor f \rfloor$  is the set of values defined by:

$$\begin{aligned} \lfloor e \rfloor &= \text{any} & \lfloor Xa \rfloor &= \text{any} \\ \lfloor p \Rightarrow f \rfloor &= \lfloor p \rfloor \& \lfloor f \rfloor & \lfloor \mu X.f \rfloor &= \lfloor f \rfloor \\ \lfloor f_1 \mid f_2 \rfloor &= \lfloor f_1 \rfloor \mid \lfloor f_2 \rfloor & \lfloor \text{groupby } f \rfloor &= [\text{any}^*] \\ \lfloor (f_1, f_2) \rfloor &= (\lfloor f_1 \rfloor, \lfloor f_2 \rfloor) & \lfloor \text{orderby } f \rfloor &= [\text{any}^*] \\ \lfloor f_1; f_2 \rfloor &= \lfloor f_1 \rfloor \\ \lfloor \{ \ell_1 : f_1, \dots, \ell_n : f_n, \dots \} \rfloor &= \{ \ell_1 : \lfloor f_1 \rfloor, \dots, \ell_n : \lfloor f_n \rfloor, \dots \} \end{aligned}$$

It is easy to show that an argument included in the accepted type is a necessary (but not sufficient, because of the cases for composition and recursion) condition for the evaluation of a filter not to fail:

**Lemma 9.** Let  $f$  be a filter and  $v$  be a value such that  $v \notin \lfloor f \rfloor$ . For every  $\gamma, \delta$ , if  $\delta; \gamma \vdash_{\text{eval}} f(v) \rightsquigarrow r$ , then  $r \equiv \Omega$ .

The proof is a straightforward induction on the structure of the derivation, and is detailed in Appendix C. The last two auxiliary definitions we need are related to product and record types. In the presence of unions, the most general form for a product type is a finite union of products (since intersections distribute on products). For instance consider the type

$$(\text{int}, \text{int}) \mid (\text{string}, \text{string})$$

This type denotes the set of pairs for which either both projections are `int` or both projections are `string`. A type such as

$$(\text{int} \mid \text{string}, \text{int} \mid \text{string})$$

is less precise, since it also allows pairs whose first projection is an `int` and second projection is a `string` and *vice versa*. We see that it is necessary to manipulate finite unions of products (and similarly for records), and therefore, we introduce the following notations:

**Lemma 10 (Product decomposition).** Let  $t \in \mathbf{Types}$  such that  $t \leq (\text{any}, \text{any})$ . A product decomposition of  $t$ , denoted by  $\pi(t)$  is a set of types:

$$\pi(t) = \{(t_1^1, t_2^1), \dots, (t_1^n, t_2^n)\}$$

such that  $t = \bigvee_{t_i \in \pi(t)} t_i$ . For a given product decomposition, we say that  $n$  is the rank of  $t$ , noted  $\text{rank}(t)$ , and use the notation  $\pi_i^j(t)$  for the type  $t_i^j$ .

There exist several suitable decompositions whose details are out of the scope of this paper. We refer the interested reader to [15] and [23] for practical algorithms that compute such decompositions for any subtype of  $(\text{any}, \text{any})$  or of  $\{..\}$ . These notions of decomposition, rank and projection can be generalized to records:

**Lemma 11 (Record decomposition).** Let  $t \in \mathbf{Types}$  such that  $t \leq \{..\}$ . A record decomposition of  $t$ , denoted by  $\rho(t)$  is a finite set of types  $\rho(t) = \{r_1, \dots, r_n\}$  where each  $r_i$  is either of the form  $\{\ell_1^i : t_1^i, \dots, \ell_{n_i}^i : t_{n_i}^i\}$  or of the form  $\{\ell_1^i : t_1^i, \dots, \ell_{n_i}^i : t_{n_i}^i, ..\}$  and such that  $t = \bigvee_{r_i \in \rho(t)} r_i$ . For a given record decomposition, we say that  $n$  is the rank of  $t$ , noted  $\text{rank}(t)$ , and use the notation  $\rho_\ell^j(t)$  for the type of label  $\ell$  in the  $j^{\text{th}}$  component of  $\rho(t)$ .

In our calculus we have three different sets of variables. The set **Vars** of term variables, ranged over by  $x, y, \dots$ , introduced in patterns and used in expressions and in arguments of calls of recursive filters. The set **RVars** of term recursion variables, ranged over by  $X, Y, \dots$  and that are used to define recursive filters. The set **TVars** of type recursion variables, ranged over by  $T, U, \dots$  used

to define recursive types. In order to use them we need to define three different environments:  $\Gamma : \mathbf{Vars} \rightarrow \mathbf{Types}$  denoting *type environments* that associate term variables with their types;  $\Delta : \mathbf{RVars} \rightarrow \mathbf{Filters}$  denoting *definition environments* that associate each filter recursion variable with the body of its definition;  $M : \mathbf{RVars} \times \mathbf{Types} \rightarrow \mathbf{TVars}$  denoting *memoization environments* which record that the call of a given recursive filter on a given type yielded the introduction of a fresh recursion type variable. Our typing rules, thus work on judgments of the form  $\Gamma; \Delta; M \vdash f(t) : t'$  stating that applying  $f$  to an expression of type  $t$  in the environments  $\Gamma, \Delta, M$  yields a result of type  $t'$ . This judgment can be derived with the set of rules given in Figure 2.

These rules are straightforward, when put side by side with the dynamic semantics of filters, given in Section 3. It is clear that this type system simulates *at the level of types* the computations that are carried out by filters on values at runtime. For instance, rule [FIL-EXPR] calls the typing function of the host language to determine the type of an expression  $e$ . Rule [FIL-PROD] applies a product filter recursively on the first and second projection for each member of the product decomposition of the input type and returns the union of all result types. Rule [FIL-REC] for records is similar, recursively applying sub-filters label-wise for each member of the record decomposition and returning the union of the resulting record types. As for the pattern filter (rule [FIL-PAT]), its subfilter  $f$  is typed in the environment augmented by the mapping  $t/p$  of the input type against the pattern (*cf.* Theorem 5). The typing rule for the union filter, [FIL-UNION] reflects the first match policy: when typing the second branch, we know that the first was not taken, hence that at runtime the filtered value will have a type that is in  $t$  but not in  $\lfloor f_1 \rfloor$ . Notice that this is *not* ensured by the definition of accepted type — which is a rough approximation that discards grosser errors but, as we stressed right after its definition, is not sufficient to ensure that evaluation of  $f_1$  will not fail — but by the type system itself: the premises *check* that  $f_1(t_1)$  is well-typed which, by induction, implies that  $f_1$  will never fail on values of type  $t_1$  and, ergo, that these values will never reach  $f_2$ . Also, we discard from the output type the contribution of the branches that cannot be taken, that is, branches whose accepted type have an empty intersection with the input type  $t$ . Composition (rule [FIL-COMP]) is straightforward. In this rule, the restriction that  $f_1$  is a filter with no open recursion variable ensures that its output type  $s$  is also a type without free recursion variables and, therefore, that we can use it as input type for  $f_2$ . The next three rules work together. The first, [FIL-FIX] introduces for a recursive filter a fresh recursion variable for its output type. It also memoize in  $\Delta$  that the recursive filter  $X$  is associated with a body  $f$  and in  $M$  that for an input filter  $X$  and an input type  $t$ , the output type is the newly introduced recursive type variable. When dealing with a recursive call  $X$  two situations may arise. One possibility is that it is the first time the filter  $X$  is applied to the input type  $t$ . We therefore introduce a fresh type variable  $T$  and recurse, replacing  $X$  by its definition  $f$ . Otherwise, if the input type has already been encountered while typing the filter variable  $X$ , we can return its memoized type, a type variable  $T$ . Finally, Rule [FIL-ORDBY] and Rule [FIL-GRPBYP] handle the special cases of `groupby` and `orderby` filters. Their typing is explained in the following section.

#### 4.4 Typing of orderby and groupby

While the “structural” filters enjoy simple, compositional typing rules, the ad-hoc operations `orderby` and `groupby` need specially crafted rules. Indeed it is well known that when transformation languages have the ability to compare data values type-checking (and also type inference) becomes undecidable (*eg.* see [1, 2]). We therefore provide two typing approximations that yield a good



$\frac{[\text{FIL-EXPR}]}{\Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} e(t) : \text{type}(\Gamma, e, )}$	$\frac{[\text{FIL-PAT}]}{\Gamma \cup t/p \ ; \Delta \ ; M \vdash_{\text{fil}} f(t) : s} \quad t \leq \wr p \wr \& \wr f \wr$	$\frac{[\text{FIL-PROD}]}{i=1..\text{rank}(t), \ j=1, 2 \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f_j(\pi_j^i(t)) : s_j^i} \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} (f_1, f_2)(t) : \bigvee_{i=1..\text{rank}(t)} (s_1^i, s_2^i)$
$\frac{[\text{FIL-REC}]}{i=1..\text{rank}(t), \ j=1..m \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f_j(\rho_{\ell_j}^i(t)) : s_j^i} \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} \{\ell_1; f_1, \dots, \ell_m; f_m, ..\}(t) : \bigvee_{i=1..\text{rank}(t)} \{\ell_1 s_1^i, \dots, \ell_m s_m^i, ..\}$	$\frac{[\text{FIL-UNION}]}{i=1, 2 \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f_i(t_i) : s_i} \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f_1   f_2(t) : \bigvee_{\{i   s_i \neq \text{empty}\}} s_i \quad \begin{array}{l} t \leq \wr f_1 \wr \wr f_2 \wr \\ t_1 = t \& \wr f_1 \wr \\ t_2 = t \& \neg \wr f_1 \wr \end{array}$	
$\frac{[\text{FIL-COMP}]}{\Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f_1(t) : s \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f_2(s) : s'} \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f_1; f_2(t) : s'$	$\frac{[\text{FIL-FIX}]}{\Gamma \ ; \Delta, (X \mapsto f) \ ; M, ((X, t) \mapsto T) \vdash_{\text{fil}} f(t) : s} \quad T \text{ fresh} \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} (\mu X. f)(t) : \mu T. s$	
$\frac{[\text{FIL-CALL-NEW}]}{\Gamma \ ; \Delta \ ; M, ((X, t) \mapsto T) \vdash_{\text{fil}} \Delta(X)(t) : t'} \quad \begin{array}{l} t = \text{type}(\Gamma, a) \\ (X, t) \notin \text{dom}(M) \\ T \text{ fresh} \end{array} \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} (Xa)(s) : \mu T. t'$	$\frac{[\text{FIL-CALL-MEM}]}{\Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} (Xa)(s) : M(X, t)} \quad \begin{array}{l} t = \text{type}(\Gamma, a) \\ (X, t) \in \text{dom}(M) \end{array}$	
$\frac{[\text{FIL-ORDBY}]}{\forall t_i \in \text{item}(t) \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f(t_i) : s_i} \quad \begin{array}{l} t \leq [\text{any*}] \\ \bigvee s_i \text{ is ordered} \end{array} \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} (\text{orderby } f)(t) : \text{OrderBy}(t)$	$\frac{[\text{FIL-GRPB Y}]}{\forall t_i \in \text{item}(t) \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} f(t_i) : s_i} \quad t \leq [\text{any*}] \quad \Gamma \ ; \Delta \ ; M \vdash_{\text{fil}} (\text{groupby } f)(t) : [(\bigvee_i s_i, \text{OrderBy}(t)) *]$	

**Figure 2.** Type inference algorithm for filter application

compromise between precision and decidability. First we define an auxiliary function over sequence types:

**Definition 12 (Item set).** Let  $t \in \mathbf{Types}$  such that  $t \leq [\text{any*}]$ . The *item set* of  $t$  denoted by  $\text{item}(t)$  is defined by:

$$\begin{aligned} \text{item}(\text{empty}) &= \emptyset \\ \text{item}(t) &= \text{item}(t \& (\text{any}, \text{any})) \quad \text{if } t \not\leq (\text{any}, \text{any}) \\ \text{item}\left(\bigvee_{1 \leq i \leq \text{rank}(t)} (t_i^1, t_i^2)\right) &= \bigcup_{1 \leq i \leq \text{rank}(t)} (\{t_i^1\} \cup \text{item}(t_i^2)) \end{aligned}$$

The first and second line in the definition ensure that  $\text{item}()$  returns the empty set for sequence types that are not products, namely for the empty sequence. The third line handles the case of non-empty sequence type. In this case  $t$  is a finite union of products, whose first components are the types of the “head” of the sequence and second components are recursively the types of the tails. Note also that this definition is well-founded. Since types are regular trees the number of distinct types accumulated by  $\text{item}()$  is finite. We can now defined typing rules for the `orderby` and `groupby` operators.

**orderby  $f$ :** The `orderby` filter uses its argument filter  $f$  to compute a key from each element of the input sequence and then returns the same sequence of elements, sorted with respect to their key. Therefore, while the types of the elements in the result are still known, their order is lost. We use  $\text{item}()$  to compute the output type of an `orderby` application:

$$\text{OrderBy}(t) = [(\bigvee_{t_i \in \text{item}(t)} t_i) *]$$

**groupby  $f$ :** The typing of `orderby` can be used to give a rough approximation of the typing of `groupby` as stated by rule [FIL-GRPB Y]. In words, we obtain a list of pairs where the key component is the result type of  $f$  applied to the items of the sequence, and use `OrderBy` to shuffle the order of the list. A far more precise typing of `groupby` that keeps track of the relation between list elements and their images via  $f$  is given in Appendix E.

#### 4.5 Soundness, termination, and complexity

The soundness of the type inference system is given by the property of subject reduction for filter application

**Theorem 13 (subject reduction).** *If  $\emptyset; \emptyset; \emptyset \vdash_{\text{fil}} f(t) : s$ , then for all  $v : t$ ,  $\emptyset; \emptyset \vdash_{\text{eval}} f(v) \rightsquigarrow r$  implies  $r : s$ .*

whose full proof is given in Appendix C. It is easy to write a filter for which the type inference algorithm, that is the deduction of  $\vdash_{\text{fil}}$ , does not terminate:  $\mu X. x \Rightarrow X(x, x)$ . The deduction of  $\Gamma; \Delta; M \vdash_{\text{fil}} f(t) : s$  simulates an (abstract) execution of the filter  $f$  on the type  $t$ . Since filters are Turing complete, then in general it is not possible to decide whether the deduction of  $\vdash_{\text{fil}}$  for a given filter  $f$  will terminate for every input type  $t$ . For this reason we define a static analysis  $\text{Check}(f)$  for filters that ensures that if  $f$  passes the analysis, then for every input type  $t$  the deduction of  $\Gamma; \Delta; M \vdash_{\text{fil}} f(t) : s$  terminates. For space reasons the formal definition of  $\text{Check}(f)$  is relegated to Appendix A, but its behavior can be easily explained. Imagine that a recursive filter  $f$  is applied to some input type  $t$ . The algorithm tracks all the recursive calls occurring in  $f$ ; next it performs one step of reduction of each recursive call by unfolding the body; finally it checks in this unfolding that if a variable occurs in the argument of a recursive call, then it is bound to a type that is a subtree of the original type  $t$ . In other words, the analysis verifies that in the execution of the derivation for  $f(t)$  every call to  $s/p$  for some type  $s$  and pattern  $p$  always yields a type environment where variables used in recursive calls are bound to subtrees of  $t$ . This implies that the rule [FIL-CALL-NEW] will always memoize for a given  $X$ , types that are obtained from the arguments of the recursive calls of  $X$  by replacing their variables with a subtree of the original type  $t$  memoized by the rule [FIL-FIX]. Since  $t$  is regular, then it has finitely many distinct subtrees, thus [FIL-CALL-NEW] can memoize only finitely many distinct types, and therefore the algorithm terminates.

More precisely, the analysis proceeds in two passes. In the first pass the algorithm tracks all recursive filters and for each of them it (i) marks the variables that occur in the arguments of its recursive calls, (ii) assigns to each variable an abstract identifier representing the subtree of the input type to which the variable will be bound at the initial call of the filter, and (iii) it returns the set of all types obtained by replacing variables by the associated abstract identifier in each argument of a recursive call. The last set intuitively represents all the possible ways in which recursive calls can shuffle and recompose the subtrees forming the initial input type. The second phase of the analysis first abstractly reduces by one step each recursive filter by applying it on the set of types collected in the first phase of the analysis and then checks whether, after this reduction, all the variables marked in the first phase (ie, those that occur in arguments of recursive calls) are still bound to subtrees of the initial input type: if this checks fails, then the filter is rejected.

It is not difficult to see that the type inference algorithm converges if and only if for every input type there exists a integer  $n$  such that after  $n$  recursive calls the marked variables are bound only to subtrees of the initial input type (or to something that does not depend on it, of course). Since deciding whether such an  $n$  exists is not possible, our analysis checks whether for all possible input types a filter satisfies it for  $n=1$ , that is to say, that at every recursive call its marked variables satisfy the property; otherwise it rejects the filter.

**Theorem 14 (Termination).** *If  $\text{Check}(f)$ , then for every type  $t$  the deduction of  $\Gamma ; \emptyset ; \emptyset \vdash_{\text{fi}} f(t) : s$  is in 2-EXPTIME. Furthermore, if  $t$  is given as a non-deterministic tree automaton (NTA) then  $\Gamma ; \emptyset ; \emptyset \vdash_{\text{fi}} f(t) : s$  is in EXPTIME, where the size of the problem is  $|f| \times |t|$ .*

(for proofs see Appendix A for termination and Appendix D for complexity). This complexity result is in line with those of similar formalisms. For instance in [20], it is shown that type-checking non deterministic top-down tree transducers is in EXPTIME when the input and output types are given by a NTA.

All filters defined in this paper (excepted those in Appendix B) pass the analysis. As an example consider the filter `rotate` that applied to a list returns the same list with the first element moved to the last position (and the empty list if applied to the empty list):

$$\mu X. ((x, (y, z)) \Rightarrow (y, X(x, z)) \mid w \Rightarrow w)$$

The analysis succeeds on this filter. If we denote by  $\iota_x$  the abstract subtree bound to the variable  $x$ , then the recursive call will be executed on the abstract argument  $(\iota_x, \iota_z)$ . So in the unfolding of the recursive call  $x$  is bound to  $\iota_x$ , whereas  $y$  and  $z$  are bound to two distinct subtrees of  $\iota_z$ . The variables in the recursive call,  $x$  and  $z$ , are thus bound to subtrees of the original tree (even though the argument of the recursive call is *not* a subtree of the original tree), therefore the filter is accepted. In order to appreciate the precision of the inference algorithm consider the type  $[\text{int}^+ \text{ bool}^+]$ , that is, the type of lists formed by some integers (at least one) followed by some booleans (at least one). For the application of `rotate` to an argument of this type our algorithm *statically* infers the most precise type, that is,  $[\text{int}^* \text{ bool}^+ \text{ int} \text{ int}] \mid [\text{bool}^* \text{ int} \text{ bool}]$ .

Generic filters are Turing complete. However, requiring that  $\text{Check}()$  holds —meaning that the filter is typeable by our system— restricts the expressive power of our filters by preventing them from *recomposing* a new value before doing a recursive call. For instance, it is not possible to typecheck a filter which reverses the elements of a sequence. Determining the exact class of transformations that typeable filters can express is challenging. However it is possible to show (cf. Appendix F) that typeable filters are strictly more expressive than top-down tree transducers with regular look-ahead, a formalism for tree transformations introduced in [14]. The intuition about this result can be conveyed by an example. Consider the tree:

$$a(u_1(\dots(u_n()))v_1(\dots(v_m())))$$

that is, a tree whose root is labeled  $a$  with two children, each being a monadic tree of height  $n$  and  $m$ , respectively. Then it is not possible to write a top-down tree transducer with regular look-ahead that creates the tree

$$a(u_1(\dots(u_n(v_1(\dots v_m())))))$$

which is just the concatenation of the two children of the root, seen as sequences, a transformation that can be easily programmed by typeable filters. The key difference in expressive power comes from the fact that filters are evaluated with an *environment* that binds capture variables to sub-trees of the input. This feature is essential to encode sequence concatenation and sequence flattening —two pervasive operations when dealing with sequences— that cannot be expressed by top-down tree transducers with regular look-ahead.

## 5. Jaql

In this Section, we show how filters can be used to capture some popular languages for processing data on the Cloud. We consider Jaql [18], a query language for JSON developed by IBM. We give translation rules from a subset of Jaql into filters.

**Definition 15 (Jaql expressions).** We use the following simplified grammar for Jaql (where we distinguish simple expressions, ranged over by  $e$ , from “core expressions” ranged over by  $k$ ).

$e ::= c$	(constants)
$  x$	(variables)
$  \$$	(current value)
$  [e, \dots, e]$	(arrays)
$  \{ e:e, \dots, e:e \}$	(records)
$  e.l$	(field access)
$  \text{op}(e, \dots, e)$	(function call)
$  e \rightarrow k$	(pipe)
$k ::= \text{filter } (\text{each } x) ? e$	(filter)
$  \text{transform } (\text{each } x) ? e$	(transform)
$  \text{expand } ((\text{each } x) ? e) ?$	(expand)
$  \text{group } ((\text{each } x) ? \text{ by } x=e (\text{as } x) ?) ? \text{ into } e$	(grouping)

### 5.1 Built-in filters

In order to ease the presentation we extend our syntax by adding “filter definitions” (already informally used in the introduction) to filters and “filter calls” to expressions:

$$e ::= \text{let filter } F[F_1, \dots, F_n] = f \text{ in } e \quad (\text{filter defn.})$$

$$f ::= F[f, \dots, f] \quad (\text{call})$$

where  $F$  ranges over *filter names*. The mapping for most of the language we consider rely on the following built-in filters.

```
let filter Filter [F] =  $\mu X$ .
  'nil  $\Rightarrow$  'nil
  | ((x, xs), tl)  $\Rightarrow$  (X(x, xs), X(tl))
  | (x, tl)  $\Rightarrow$  Fx ; ('true  $\Rightarrow$  (x, X(tl))) | 'false  $\Rightarrow$  X(tl)

let filter Transform [F] =  $\mu X$ .
  'nil  $\Rightarrow$  'nil
  | ((x, xs), tl)  $\Rightarrow$  (X(x, xs), X(tl))
  | (x, tl)  $\Rightarrow$  (Fx, X(tl))

let filter Expand =  $\mu X$ .
  'nil  $\Rightarrow$  'nil
  | ('nil, tl)  $\Rightarrow$  X(tl)
  | ((x, xs), tl)  $\Rightarrow$  (x, X(xs, tl))
```

### 5.2 Mapping

Jaql expressions are mapped to our expressions as follows (where  $\$$  is a distinguished expression variable interpreting Jaql’s  $\$$ ):

$\llbracket c \rrbracket$	$= c$
$\llbracket x \rrbracket$	$= x$
$\llbracket \$ \rrbracket$	$= \$$
$\llbracket [e_1:e'_1, \dots, e_n:e'_n] \rrbracket$	$= \{ \llbracket e_1 \rrbracket : \llbracket e'_1 \rrbracket, \dots, \llbracket e_n \rrbracket : \llbracket e'_n \rrbracket \}$
$\llbracket e.l \rrbracket$	$= \llbracket e \rrbracket.l$
$\llbracket \text{op}(e_1, \dots, e_n) \rrbracket$	$= \text{op}(\llbracket e_1 \rrbracket, \dots, \llbracket e_n \rrbracket)$
$\llbracket [e_1, \dots, e_n] \rrbracket$	$= (\llbracket e_1 \rrbracket, \dots, \llbracket e_n \rrbracket, \text{'nil'})$
$\llbracket e \rightarrow k \rrbracket$	$= \llbracket e \rrbracket ; \llbracket k \rrbracket_F$

Jaql core expressions are mapped to filters as follows:

$$\begin{aligned}
\llbracket \text{filter } e \rrbracket_F &= \llbracket \text{filter each } \$e \rrbracket_F \\
\llbracket \text{filter each } x \ e \rrbracket_F &= \text{Filter } [x \Rightarrow \llbracket e \rrbracket] \\
\llbracket \text{transform } e \rrbracket_F &= \llbracket \text{transform each } \$e \rrbracket_F \\
\llbracket \text{transform each } x \ e \rrbracket_F &= \text{Transform } [x \Rightarrow \llbracket e \rrbracket] \\
\llbracket \text{expand each } x \ e \rrbracket_F &= \llbracket \text{expand} \rrbracket_F ; \llbracket \text{transform each } x \ e \rrbracket_F \\
\llbracket \text{expand} \rrbracket_F &= \text{Expand} \\
\llbracket \text{group into } e \rrbracket_F &= \llbracket \text{group by } y=\text{true into } e \rrbracket_F \\
\llbracket \text{group by } y=e_1 \text{ into } e_2 \rrbracket_F &= \llbracket \text{group each } \$ \text{ by } y=e_1 \text{ into } e_2 \rrbracket_F \\
\llbracket \text{group each } x \text{ by } y=e_1 \text{ into } e_2 \rrbracket_F &= \\
&\quad \llbracket \text{group each } x \text{ by } y = e_1 \text{ as } \$ \text{ into } e_2 \rrbracket_F \\
\llbracket \text{group each } x \text{ by } y = e_1 \text{ as } g \text{ into } e_2 \rrbracket_F &= \\
&\quad \text{groupby } x \Rightarrow [e_1] ; \text{Transform } [(y,g) \Rightarrow \llbracket e_2 \rrbracket]
\end{aligned}$$

This translation defines the (first, in our knowledge) formal semantics of Jaql. Such a translation is *all* that is needed to define the semantics of a NoSQL language and, as a bonus, endow it with the type inference system we described *without requiring any modification of the original language*. No further action is demanded since the machinery to exploit it is all developed in this work.

As for typing, every Jaql expression is encoded into a filter for which type-checking is ensured to terminate: *Check()* holds for *Filter*[], *Transform*[], and *Expand* (provided it holds also for their arguments) since they only perform recursive calls on recombinations of subtrees of their input; by its definition, the encoding does not introduce any new recursion and, hence, it always yields a composition and application of filters for which *Check()* holds.

### 5.3 Examples

To show how we use the encoding, let us encode the example of the introduction. For the sake of the concision we will use filter definitions (rather than expanding them in details). We use *Fil* and *Sel* defined in the introduction, *Expand* and *Transform*[] defined at the beginning of the section, the encoding of Jaql's field selection as defined in Section 3.3, and finally *Head* that returns the first element of a sequence and a family of recursive filters *Rgrpi* with  $i \in \mathbb{N}^+$  both defined below:

```

let filter Head = 'nil => null | (x,xs) => x
let filter Rgrpi = 'nil => 'nil
                  | ((i,x),tail) => (x , Rgrpi tail)
                  | _ => Rgrpi tail

```

Then, the query in the introduction is encoded as follows

```

1  [employees depts];
2  [Sel Fil];
3  [Transform[x =>(1,x)] Transform[x =>(2,x)]];
4  Expand;
5  groupby ( (1,$)=>$.dept | (2,$)=>$.depid );
6  Transform[(g,1)=>(
7      [(1; Rgrp1) (1; Rgrp2)];
8      [es ds] =>
9          { dept: g,
10             deptName: (ds ; Head).name),
11             numEmps: count(es) } )])

```

In words, we perform the selection on employees and filter the departments (lines 1-2); we tag each element by 1 if it comes from employees, and by 2 if it comes from departments (line 3); we merge the two collections (line 4); we group the heterogeneous list according to the corresponding key (line 5); then for each element of the result of grouping we capture in *g* the key (line 6), split the group into employees and depts (line 7), capture each subgroup into the corresponding variable (*ie*, *es* and *ds*) (line 8) and return the expression specified in the query after the “*into*” (lines 8-10). The general definition of the encoding for the co-grouping is given in Appendix H.

Let us now illustrate how the above composition of filters is typed. Consider an instance where:

- *employees* has type  $[ \text{Remp}^* ]$ , where  $\text{Remp} \equiv \{ \text{dept: int, income: int, ...} \}$
- *depts* has type  $[ (\text{Rdep} \mid \text{Rbranch})^* ]$ , where  $\text{Rdep} \equiv \{ \text{depid: int, name: string, size: int} \}$   $\text{Rbranch} \equiv \{ \text{brid: int, name: string} \}$  (this type is a subtype of *Dept* as defined in the introduction)

The global input type is therefore (line 1)

$$[ [ \text{Remp}^* ] [ (\text{Rdep} \mid \text{Rbranch})^* ] ]$$

which becomes, after selection and filtering (line 2)

$$[ [ \text{Remp}^* ] [ \text{Rdep}^* ] ]$$

(note how all occurrences of *Rbranch* are ignored by *Fil*). Tagging with an integer (line 3) and flattening (line 4) yields

$$[ (1, \text{Remp})^* (2, \text{Rdep})^* ]$$

which illustrates the precise typing of products coupled with singleton types (*ie*, 1 instead of *int*). While the *groupby* (line 5) introduces an approximation the dependency between the tag and the corresponding type is kept

$$[ (\text{int}, [ ((1, \text{Remp}) \mid (2, \text{Rdep}))^* ]) ]$$

Lastly the transform is typed exactly, yielding the final type

$$[ \{ \text{dept: int, deptName: string} \mid \text{null, numEmps: int} \}^* ]$$

Note how *null* is retained in the output type (since there may be employees without a department, then *Head* may be applied to an empty list returning *null*, and the selection of name of *null* returns *null*). For instance suppose to pipe the Jaql grouping defined in the introduction into the following Jaql expression, in order to produce a printable representation of the records of the result

```

transform each x (
  (x.deptName)@"@"(to_string x.dep)@"@"(x.numEmps))

```

where *@* denotes string concatenation and *to\_string* is a conversion operator (from any type to string). The composition is ill-typed for three reasons: the field *dept* is misspelled as *dep*, *x.numEmps* is of type *int* (so it must be applied to *to\_string* before concatenation), and the programmer did not account for the fact that the value stored in the field *deptName* may be *null*. The encoding produces the following lines to be appended to the previous code:

```

12 Transform[ x =>
13   (x.deptName)@"@"(to_string x.dep)@"@"(x.numEmps)]

```

in which all the three errors are detected by our type system. A subtler example of error is given by the following alternative code

```

12 Transform[
13   { dept : d, deptName: n&String, numEmps: e } =>
14     n @ ":" @ (to_string d) @ ":" @ (to_string e)
15   | { deptName: null, .. } => ""
16   | _ => "Invalid department" ]

```

which corrects all the previous errors but adds a new one since, as detected by our type system, the last branch can be never selected. As we can see, our type-system ensures soundness, forcing the programmer to handle exceptional situations (as in the *null* example above) but is also precise enough to detect that some code paths can never be reached.

In order to focus on our contributions we kept the language of types and filters simple. However there already exists several contributions on the types and expressions used here. Two in particular are worth mentioning in this context: recursive patterns and XML.

Definition 3 defines patterns inductively but, alternatively, we can consider the (possibly infinite) regular trees *coinductively* generated by these productions and, on the lines of what is done in CDuce, use the recursive patterns so obtained to encode regular expressions patterns (see [4]). Although this does not enhance ex-



pressiveness, it greatly improves the writing of programs since it makes it possible to capture distinct subsequences of a sequence by a single match. For instance, when a sequence is matched against a pattern such as `[ (int as x | bool as y | _)* ]`, then `x` captures (the list of) all integer elements (capture variables in regular expression patterns are bound to lists), `y` captures all Boolean elements, while the remaining elements are ignored. By such patterns, co-grouping can be encoded without the `Rgrp`. For instance, the transform in lines 6-11 can be more compactly rendered as:

```
6 Transform[(g,[ ((1,es)|(2,ds))* ]) =>
7     { dept: g,
8       deptName: (ds;Head).name,
9       numEmps: count(es) }]
```

For what concerns XML, the types used here were originally defined for XML, so it comes as a no surprise that they can seamlessly express XML types and values. For example `CDuce` uses the very same types used here to encode both XML types and elements as triples, the first element being the tag, the second a record representing attributes, and the third a heterogeneous sequence for the content of the element. Furthermore, we can adapt the results of [9] to encode forward XPath queries in filters. Therefore, it requires little effort to use the filters presented here to encode languages such as JSONiq [12] designed to integrate JSON and XML, or to precisely type regular expressions, the import/export of XML data, or XPath queries embedded in Jaql programs. This is shown in the section that follows.

## 6. JSON, XML, Regex

There exist various attempts to integrate JSON and XML. For instance JSONiq [12] is a query language designed to allow XML and JSON to be used in the same query. The motivation is that JSON and XML are both widely used for data interchange on the Internet. In many applications, JSON is replacing XML in Web Service APIs and data feeds, while more and more applications support both formats. More precisely, JSONiq embeds JSON into an XML query language (XQuery), but it does it in a stratified way: JSONiq does not allow XML nodes to contain JSON objects and arrays. The result is thus similar to OCamlDuce, the embedding of `CDuce`'s XML types and expressions into OCaml, with the same drawbacks.

Our type system is derived from the type system of `CDuce`, whereas the theory of filters was originally designed to use `CDuce` as an host language. As a consequence XML types and expressions can be seamlessly integrated in the work presented here, without any particular restriction. To that end it suffices to use for XML elements and types the same encoding used in the implementation of `CDuce`, where an XML element is just a triple formed by a tag (here, an expression), a record (whose labels are the attributes of the element), and a sequence (of characters and or other XML elements) denoting its content. So for instance the following element

```
<product system="US-size">
  <number>557</number>
  <name>Blouse</name>
</product>
```

is encoded by the following triple:

```
("product" , { system : "US-size" } ,
 [
   ("number" , {} , [ 557 ])
   ("name" , {}, "Blouse")
 ]
)
```

and this latter, with the syntactic sugar defined for `CDuce`, can be written as:

```
<product system="US-size">[
  <number>[ 557 ]
  <name>[ Blouse ]
]
```

Clearly in our system there are no restrictions in merging and nesting JSON and XML and no further extension is required to our system to define XML query and processing expressions. Just, the introduction of syntactic sugar to make the expressions readable, seems helpful:

$$\begin{aligned} e &::= \langle e \ e \rangle e \\ f &::= \langle f \ f \rangle f \end{aligned}$$

The system we introduced here is already able to reproduce (and type) the same transformations as in JSONiq, but without the restrictions and drawbacks of the latter (this is why we argue that it is better to extend NoSQL languages with XML primitives directly derived from our system rather than to use our system to encode JSONiq). For instance, the example given in the JSONiq draft to show how to render in Xhtml the following JSON data:

```
{
  "col labels" : ["singular", "plural"],
  "row labels" : ["1p", "2p", "3p"],
  "data" :
    [
      ["spinne", "spinnen"],
      ["spinnst", "spinnt"],
      ["spinnt", "spinnen"]
    ]
}
```

can be encoded in the filters presented in this work (with the new syntactic sugar) as:

```
{ "col labels" : c1 ,
  "row labels" : r1 ,
  "data" : d1
} =>
<table border="1" cellpadding="1" cellspacing="2">[
  <tr>[ <th>[ ] !(c1; Transform[ x -> <th>x ]) ]
    !(r1; Transform[ h ->
      <tr>[ <th>h !(d1; Transform[ x -> <td>x ]) ]
    ]
  )
]
```

(where `!` expands a subsequence in the containing sequence). The resulting Xhtml document is rendered in a web browser as:

	singular	plural
1p	spinne	spinnen
2p	spinnst	spinnt
3p	spinnt	spinnen

Similarly, Jaql built-in libraries include functions to convert and manipulate XML data. So for example as it is possible in Jaql to embed SQL queries, so it is possible to evaluate XPath expressions, by the function `xpath()` which takes two arguments, an XML document and a string containing an xpath expression—eg, `xpath(read(seq(("conf/addr.xml")), "content/city")` —. Filters can encode forward XPath expressions (see [23]) and precisely type them. So while in the current implementation there is no check of the type of the result of an external query (nor for

XPath or for SQL) and the XML document is produced independently from Jaql, by encoding (forward) XPath into filters we can not only precisely type calls to Jaql’s `xpath()` function but also feed them with documents produced by Jaql expressions.

Finally, the very same regular expressions types that are used to describe heterogeneous sequences and, in particular, the content of XML elements, can be used to type regular expressions. Functions working on regular expressions (regex for short) form, in practice, yet another domain specific language that is embedded in general purpose languages in an untyped or weakly typed (typically, every result is of type `string`) way. Recursive patterns can straightforwardly encode regexp matching. Therefore, by combining the pattern filter with other filters it is possible to encode any regexp library, with the important advantage that, as stated by Theorem 4, the set of values (respectively, strings) accepted by a pattern (respectively, by a regular expression), can be precisely computed and can be expressed by a type. So a function such as the built-in Jaql’s function `regex_extract()`, which extracts the substrings that match a given regex, can be easily implemented by filters and precisely typed by our typing rules. Typing will then amount to intersect the type of the string (which can be more precise than just `string`) with the type accepted by the pattern that encodes the regex at issue.

## 7. Programming with filters

Up to now we have used the filters to encode operators hard-coded in some languages in order to type them. Of course, it is possible to embed the typing technology we introduced directly into the compilers of these language so as to obtain the flexible typing that characterizes our system. However, an important aspect of filters we have ignored so far is that they can be used directly by the programmer to define user-defined operators that are typed as precisely as the hard-coded ones. Therefore a possibility is extend existing NoSQL languages by adding to their expressions the filter application  $fe$  expression.

The next problem is to decide how far to go in the definition of filters. A complete integration, that is taking for  $f$  all the definitions given so far, is conceivable but might disrupt the execution model of the host language, since the user could then define complex iterators that do not fit map-reduce or the chosen distributed compilation policy. A good compromise could be to add to the host language only filters which have “local” effects, thus avoiding to affect the map-reduce or distributed compilation execution model. The minimal solution consists in choosing just the filters for patterns, unions, and expressions:

$$f ::= e \mid p \Rightarrow f \mid f|f$$

Adding such filters to Jaql (we use the “ $\Rightarrow$ ” arrow for patterns in order to avoid confusion with Jaql’s “ $->$ ” pipe operator) would not allow the user to define powerful operators, but their use would already dramatically improve type precision. For instance we could define the following Jaql expression

```
transform ( {a:x, ..} as y => {y.*, sum:x+x} | y => y )
```

(with the convention that a filter occurring as an expression denotes its application to the current argument  $\$$ ). With this syntax, our inference system is able to deduce that feeding this expression with an argument of type  $\{[a?:int, c:bool]*\}$  returns a result of type  $\{([a:int, c:bool, sum:int] \mid [c:bool])*\}$ . This precision comes from the capacity of our inference system to discriminate between the two branches of the filter and deduce that a `sum` field will be added only if the `a` field is present. Similarly by using pattern matching in a Jaql “`filter`” expression, we can deduce that `filter ( int=>true | _=>false )` fed with any sequence of elements always returns a (possibly empty) list of in-

tegers. An even greater precision can be obtained for grouping expressions when the generation of the key is performed by a filter that discriminates on types: the result type can keep a precise correspondence between keys and the corresponding groups. As an example consider the following (extended) Jaql grouping expression:

```
group e by({town: ("Roma"|"Pisa"), ..} => "Italia"
           |{town: "Paris", ..} => "France"
           | _ => "?")
```

if  $e$  has type  $\{[town:string, addr:string]*\}$ , then the type inferred by our system for this groupby expression is

```
[ ( ("Italia", [{town:"Roma"|"Pisa", addr:string}+])
  | ("France", [{town:"Paris", addr:string}+])
  | ("?", [{town:string/("Roma"|"Pisa"|"Paris"),
               addr:string}+])
  )*]
```

which precisely associates every key to the type of the elements it groups.

Finally, in order to allow a modular usage of filters, adding just filter application to the expression of the foreign language does not suffice: parametric filter definitions are also needed.

$$e ::= fe \mid \text{let filter } F[F_1, \dots, F_n] = f \text{ in } e$$

However, on the line of what we already said about the disruption of the execution model, recursive parametric filter definitions should be probably disallowed, since a compilation according to a map-reduce model would require to disentangle recursive calls.

## 8. Commentaries

Finally, let us explain some subtler design choices for our system.

**Filter design:** The reader may wonder whether products and record filters are really necessary since, at first sight, the filter  $(f_1, f_2)$  could be encoded as  $(x, y) \Rightarrow (f_1x, f_2y)$  and similarly for records. The point is that  $f_1x$  and  $f_2y$  are expressions —and thus their pair is a filter— only if the  $f_i$ ’s are closed (ie, without free term recursion variables). Without an explicit product filter it would not be possible to program a filter as simple as the identity map,  $\mu X. \text{'nil'} \Rightarrow \text{'nil'}(h, t) \Rightarrow (h, Xt)$  since  $Xt$  is not an expression ( $X$  is a free term recursion variable). Similarly, we need an explicit record filter to process recursively defined record types such as  $\mu X. \{ \text{head:int}, \text{tail:X} \} | \text{'nil'}$ .

Likewise, one can wonder why we put in filters only the “open” record variant that copy extra fields and not the closed one. The reason is that if we want a filter to be applied only to records with exactly the fields specified in the filter, then this can be simply obtained by a pattern matching. So the filter  $\{ \ell_1:f_1, \dots, \ell_n:f_n \}$  (ie, without the trailing “ $..$ ”) can be simply introduced as syntactic sugar for  $\{ \ell_1:\text{any}, \dots, \ell_n:\text{any} \} \Rightarrow \{ \ell_1:f_1, \dots, \ell_n:f_n, .. \}$

**Constructors:** The syntax for constructing records and pairs is exactly the same in patterns, types, expressions, and filters. The reader may wonder why we did not distinguish them by using, say,  $\times$  for product types or  $=$  instead of  $:$  in record values. This, combined with the fact that values and singletons have the same syntax, is a critical design choice that greatly reduces the confusion in these languages, since it makes it possible to have a unique representation

---

Syntactically we could write  $\mu X. \text{'nil'} \Rightarrow \text{'nil'}(h, t) \Rightarrow (h, (Xt)v)$  where  $v$  is any value, but then this would not pass type-checking since the expression  $(Xt)v$  must be typeable without knowing the  $\Delta$  environment (cf. rule [FILTER-APP] at the beginning of Section 4.3). We purposely stratified the system in order to avoid mutual recursion between filters and expressions.

for constructions that are semantically equivalent. Consider for instance the pattern  $(x, \{3, \text{'nil'}\})$ . With our syntax  $\{3, \text{'nil'}\}$  denotes both the product type of two singletons 3 and 'nil', or the value  $\{3, \text{'nil'}\}$ , or the singleton that contains this value. According to the interpretation we choose, the pattern can then be interpreted as a pattern that matches a product or a pattern that matches a value. If we had differentiated the syntax of singletons from that of values (eg,  $\{v\}$ ) and that of pairs from products, then the pattern above could have been written in five different ways. The point is that they all would match exactly the same sets of values, which is why we chose to have the same syntax for all of them.

**Record types:** The definition of records is redundant (both for types and patterns). Instead of the current definition we could have used just  $\{\ell:t\}$ ,  $\{\}$ , and  $\{\dots\}$ , since the rest can be encoded by intersections. For instance,  $\{\ell_1:t_1, \dots, \ell_n:t_n, \dots\} = \{\ell_1:t_1\} \& \dots \& \{\ell_n:t_n\} \& \{\dots\}$ . We opted to use the redundant definition for the sake of clarity. In order to type records with computed labels we distinguished two cases according to whether the type of a record label is finite or not. Although such a distinction is simple, it is not unrealistic. Labels with singleton types cover the (most common) case of records with statically fixed labels. The dynamic choice of a label from a statically known list of labels is a usage pattern seen in JavaScript when building an object which must conform to some interface based on a run-time value. Labels with infinite types cover the fairly common usage scenario in which records are used as dictionaries: we deduce for the expression computing the label the type `string`, thus forcing the programmer to insert some code that checks that the label is present before accessing it.

The rationale behind the typing of records was twofold. First and foremost, in this work we wanted to avoid type annotations at all costs (since there is not even a notion of schema for JSON records and collections —only the notion of basic type is defined— we cannot expect the Jaql programmer to put any kind of type information in the code). More sophisticated type systems, such as dependent types, would probably preclude type reconstruction: dependent types need a lot of annotations and this does not fit our requirements. Second, we wanted the type-system to be simple yet precise. Making the finite/infinite distinction increases typing precision at no cost (we do not need any extra machinery since we already have singleton types). Adding heuristics or complex analysis just to gain some precision on records would have blurred the main focus of our paper, which is not on typing records but on typing *transformations* on records. We leave such additions for future work.

**Record polymorphism:** The type-oriented reader will have noticed that we do not use row variables to type records, and nevertheless we have a high degree of polymorphism. Row variables are useful to type functions or transformations since they can keep track of record fields that are not modified by the transformation. In this setting we do not need them since we do not type transformations (*ie*, filters) but just the application of transformations (filters are not first-class terms). We have polymorphic typing via filters (see how the first example given in Section 7 keeps track of the `c` field) and therefore open records suffice.

**Record selection:** Some languages —typically, the dynamic ones such as Javascript, Ruby, Python— allow the label of a field selection to be computed by an expression. We considered the definition of a fine-grained rule to type expressions of the form  $e_1.e_2$ : whenever  $e_2$  is typed by a finite unions of strings, the rule would give a finite approximation of the type of the selection. However, such an extension would complex the definition of the type system, just to handle few interesting cases in which a finite union type can be deduced. Therefore, we preferred to omit its study and leave it for future work.

## 9. Related Work

In the (nested) relational (and SQL) context, many works have studied the integration of (nested)-relational algebra or SQL into general purpose programming languages. Among the first attempts was the integration of the relational model in Pascal [31] or in Smalltalk [11]. Also, monads or comprehensions [7, 33, 34] have been successfully used to design and implement query languages including a way to embed queries within host languages. Significant efforts have been done to equip those languages with type systems and type checking disciplines [8, 25, 26] and more recently [27] for integration and typing aspects. However, these approaches only support homogeneous sequences of records in the context of specific classes of queries (practically equivalent to a nested relational algebra or calculus), they do not account for records with computable labels, and therefore they are not easily transposable to a setting where sequences are heterogeneous, data are semi-structured, and queries are much more expressive.

While the present work is inspired and stems from previous works on the XML iterators, targeting NoSQL languages made the filter calculus presented here substantially different from the one of [9, 23] (dubbed XML filters in what follows), as well in syntax as in dynamic and static semantics. In [9] XML filters behave as some kind of top-down tree transducers, termination is enforced by heavy syntactic restrictions, and a *less* constrained use of the composition makes type inference challenging and requires sometimes cumbersome type annotations. While XML filters are allowed to operate by composition on the *result* of a recursive call (and, thus, simulate bottom-up tree transformations), the absence of explicit arguments in recursive calls makes programs understandable only to well-trained programmers. In contrast, the main focus of the current work was to make programs immediately intelligible to any functional programmer and make filters effective for the typing of sequence transformations: sequence iteration, element filtering, one-level flattening. The last two are especially difficult to write with XML filters (and require type annotations). Also, the integration of filters with record types (absent in [9] and just sketched in [23]) is novel and much needed to encode JSON transformations.

## 10. Conclusion

Our work addresses two very practical problems, namely the typing of NoSQL languages and a comprehensive definition of their semantics. These languages add to list comprehension and SQL operators the ability to work on heterogeneous data sets and are based on JSON (instead of tuples). Typing precisely each of these features using the best techniques of the literature would probably yield quite a complex type-system (mixing row polymorphism for records, parametric polymorphism, some form of dependent typing,...) and we are skeptical that this could be achieved without using any explicit type annotation. Therefore we explored the formalization of these languages from scratch, by defining a calculus and a type system. The thesis we defended is that all operations typical of current NoSQL languages, as long as they operate structurally (*ie*, without resorting on term equality or relations), amount to a combination of more basic bricks: our filters. On the structural side, the claim is that combining recursive records and pairs by unions, intersections, and negations suffices to capture all possible structuring of data, covering a palette ranging from comprehensions, to heterogeneous lists mixing typed and untyped data, through regular expressions types and XML schemas. Therefore, our calculus not only provides a simple way to give a formal semantics to, reciprocally compare, and combine operators of different NoSQL languages, but also offers a means to equip these languages, in their current definition (*ie*, without any type definition or annotation), with precise type inference.



As such we accounted for both components that, according to Landin, constitute the design of a language: operators and data structures. But while Landin considers the design of terms and types as independent activities we, on the contrary, advocate an approach in which the design of former is *driven* by the form of latter. Although other approaches are possible, we tried to convey the idea that this approach is nevertheless the only one that yields a type system whose precision, that we demonstrated all the work long, is comparable only to the precision obtained with hard-coded (as opposed to user-defined) operators. As such, our type inference yields and surpasses in precision systems using parametric polymorphism and row variables. The price to pay is that *transformations* are not first class: we do not type filters but just their applications. However, this seems an advantageous deal in the world of NoSQL languages where “selects” are never passed around (at least, not explicitly), but early error detection is critical, especially in the view of the cost of code deployment.

The result are filters, a set of untyped terms that can be easily included in a host language to complement in a typeful framework existing operators with user-defined ones. The requirements to include filters into a host language are so minimal that every modern typed programming language satisfies them. The interest resides not in the fact that we can add filter applications to any language, rather that filters can be used to define a smooth integration of calls to domain specific languages (eg, SQL, XPath, Pig, Regex) into general purpose ones (eg, Java, C#, Python, OCaml) so as both can share the same set of values and the same typing discipline. Likewise, even though filters provide an early prototyping platform for queries, they cannot currently be used as a final compilation stage for NoSQL languages: their operations rely on a Lisp-like encoding of sequences and this makes the correspondence with optimized bulk operations on lists awkward. Whether we can derive an efficient compilation from filters to map-reduce (recovering the bulk semantics of the high-level language) is a challenging question.

Future plans include practical experimentation of our technique: we intend to benchmark our type analysis against existing collections of Jaql programs, gauge the amount of code that is ill typed and verify on this how frequently the programmer adopted defensive programming to cope with the potential type errors.

## References

- [1] N. Alon, T. Milo, F. Neven, D. Suciu, and V. Vianu. XML with data values: typechecking revisited. In *Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, PODS '01, pages 138–149, New York, NY, USA, 2001. ACM.
- [2] N. Alon, T. Milo, F. Neven, D. Suciu, and V. Vianu. Typechecking XML views of relational databases. *ACM Trans. Comput. Logic*, 4:315–354, July 2003.
- [3] A. Behm, V. R. Borkar, M. J. Carey, R. Grover, C. Li, N. Onose, R. Vernica, A. Deutsch, Y. Papakonstantinou, and V. J. Tsotras. Asterix: towards a scalable, semistructured data platform for evolving-world models. *Distributed and Parallel Databases*, 29(3):185–216, 2011.
- [4] V. Benzaken, G. Castagna, and A. Frisch. CDuce: an XML-friendly general purpose language. In *ICFP '03, 8th ACM International Conference on Functional Programming*, pages 51–63, Uppsala, Sweden, 2003. ACM Press.
- [5] K. S. Beyer, V. Ercegovic, R. Gemulla, A. Balmin, M. Y. Eltabakh, C.-C. Kanne, F. Özcan, and E. J. Shekita. Jaql: A scripting language for large scale semistructured data analysis. *PVLDB*, 4(12):1272–1283, 2011.
- [6] S. Boag, D. Chamberlain, M. F. Fernández, D. Florescu, J. Robie, and J. Siméon. XQuery 1.0: An XML query language, W3C recommendation, 2007.
- [7] P. Buneman, L. Libkin, D. Suciu, V. Tannen, and L. Wong. Comprehension syntax. *SIGMOD Record*, 23(1):87–96, 1994.
- [8] P. Buneman, R. Nikhil, and R. Frankel. A Practical Functional Programming System for Databases. In *Proc. Conference on Functional Programming and Architecture*. ACM, 1981. 1981.
- [9] G. Castagna and K. Nguyễn. Typed iterators for XML. In *Proceeding of the 13th ACM SIGPLAN international conference on Functional programming*, ICFP '08, pages 15–26, New York, NY, USA, 2008. ACM.
- [10] H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, C. Löding, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications. <http://www.grappa.univ-lille3.fr/tata>, 2007.
- [11] G. Copeland and D. Maier. Making smalltalk a database system. In *In Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 316–325, 1984.
- [12] J. R. (editor). Jsoniq. <http://jsoniq.org>.
- [13] J. Engelfriet. Bottom-up and Top-down Tree Transformations – A comparison. *Theory of Computing Systems*, 9(2):198–231, 1975.
- [14] J. Engelfriet. Top-down tree transducers with regular look-ahead. *Mathematical Systems Theory*, 10(1):289–303, Dec. 1976.
- [15] A. Frisch. *Théorie, conception et réalisation d'un langage de programmation fonctionnel adapté à XML*. PhD thesis, Université Paris 7 Denis Diderot, 2004.
- [16] A. Frisch, G. Castagna, and V. Benzaken. Semantic subtyping: Dealing set-theoretically with function, union, intersection, and negation types. *Journal of the ACM*, 55(4):1–64, 2008.
- [17] H. Hosoya. *Foundations of XML Processing: The Tree Automata Approach*. Cambridge University Press, 2010.
- [18] Jaql. <http://code.google.com/p/jaql>.
- [19] Javascript object notation (json). <http://json.org/>.
- [20] W. Martens and F. Neven. Typechecking top-down uniform unranked tree transducers. In *Proceedings of the 9th International Conference on Database Theory*, ICDT '03, pages 64–78, London, UK, 2002. Springer-Verlag.
- [21] E. Meijer. The world according to LINQ. *ACM Queue*, 9(8):60, 2011.
- [22] E. Meijer and G. Bierman. A co-relational model of data for large shared data banks. *Communications of the ACM*, 54(4):49–58, 2011.
- [23] K. Nguyễn. *Language of Combinators for XML: Conception, Typing, Implementation*. PhD thesis, Université Paris-Sud 11, <http://www.lri.fr/~kn/files>, 2008.
- [24] Odata. <http://www.odata.org/>.
- [25] A. Ohori and P. Buneman. Type Inference in a Database Programming Language. In *ACM Conference on LISP and Functional Programming*, pages 174–183, 1988.
- [26] A. Ohori, P. Buneman, and V. Tannen. Database Programming in Machiavelli -a Polymorphic Language with Static Type Inference. *Proc. ACM SIGMOD Conference*, pages 46–57, 1989.
- [27] A. Ohori and K. Ueno. Making standard ML a practical database programming language. In M. M. T. Chakravarty, Z. Hu, and O. Danvy, editors, *ICFP*, pages 307–319. ACM, 2011.
- [28] C. Olston, B. Reed, U. Srivastava, R. Kumar, and A. Tomkins. Pig latin: a not-so-foreign language for data processing. In *SIGMOD Conference*, pages 1099–1110, 2008.
- [29] F. Özcan, D. Hoa, K. S. Beyer, A. Balmin, C. J. Liu, and Y. Li. Emerging trends in the enterprise data analytics: connecting Hadoop and DB2 warehouse. In *Proceedings of the 2011 international conference on Management of data*, SIGMOD '11, pages 1161–1164, 2011.
- [30] A. Sabry and P. Wadler. A reflection on call-by-value. In *Proceedings of the first ACM SIGPLAN international conference on Functional*

---

Only filter-encoded operators are not first class: if the host language provides, say, higher-order functions, then they stay higher-order and are typed by embedding the host type system, if any, via “foreign type calls”.

- programming*, ICFP '96, pages 13–24. ACM, 1996.
- [31] J. Schmidt and M. Mall. Pascal/R Report. Technical Report 66, Fachbereich Informatik, universit  de Hamburg, 1980.
  - [32] Squeryl: A Scala ORM and DSL for talking with Databases with minimum verbosity and maximum type safety .  
<http://squeryl.org/>.
  - [33] V. Tannen, P. Buneman, and L. Wong. Naturally embedded query languages. In *ICDT*, pages 140–154, 1992.
  - [34] P. Trinder and P. Wadler. Improving list comprehension database queries. In *Fourth IEEE Region 10 Conference (TENCON)*, pages 186–192, Nov. 1989.
  - [35] Unql. <http://www.unqlspec.org/>.

# Appendix



## A. Termination analysis algorithm

In order to deduce the result type of the application of a filter to an expression, the type inference algorithm abstractly executes the filter on the type of the expression. As explained in Section 4, the algorithm essentially analyzes what may happen to the original input type (and to its subtrees) after a recursive call and checks that, in all possible cases, every subsequent recursive call will only be applied to subtrees of the original input type. In order to track the subtrees of the original input type, we use an infinite set  $\mathbf{Ids} = \{\iota_1, \iota_2, \dots\}$  of *variable identifiers* (ranged over by, possibly indexed,  $\iota$ ). These identifiers are used to identify the subtrees of the original input type that are bound to some variable after a recursive call. Consider for instance the recursive filter:

$$\mu X.(x_1.(x_2, x_3)) \Rightarrow X(x_1, x_3) \mid \_ \Rightarrow \text{'nil'}$$

The algorithm records that at the first call of this filter each variable  $x_i$  is bound to a subtree  $\iota_i$  of the input type. The recursive call  $X(x_1, x_3)$  is thus applied to the “abstract argument”  $(\iota_1, \iota_3)$ . If we perform the substitutions for this call, then we see that  $x_1$  will be bound to  $\iota_1$  and that  $x_2$  and  $x_3$  will be bound to two subtrees of the  $\iota_3$  subtree. We do not *care* about  $x_2$  since it is not used in any recursive call. What is important is that both  $x_1$  and  $x_3$  are bound to subtrees of the original input type (respectively, to  $\iota_1$  and to a subtree of  $\iota_3$ ) and therefore, for this filter, the type inference algorithm will terminate for every possible input type.

The previous example introduces the two important concepts that are still missing for the definition of our analysis:

1. Recursive calls are applied to *symbolic arguments*, such as  $(\iota_1, \iota_3)$ , that are obtained from arguments by replacing variables by variable identifiers. Symbolic arguments are ranged over by  $A$  and are formally defined as follows:

<b>Symb Args</b>	$A ::= \iota$	<b>(variable identifiers)</b>
	$c$	<b>(constants)</b>
	$(A, A)$	<b>(pairs)</b>
	$\{\ell:A, \dots, \ell:A\}$	<b>(records)</b>
	$\perp$	<b>(indeterminate)</b>

2. We said that we “care” for some variables and disregard others. When we analyze a recursive filter  $\mu X.f$  the variables we care about are those that occur in arguments of recursive calls of  $X$ . Given a filter  $f$  and a filter recursion variable  $X$  the set of these variables is formally defined as follows

$$\text{Mark}_X(f) = \{x \mid Xa \sqsubseteq f \text{ and } x \in \text{vars}(a)\}$$

where  $\sqsubseteq$  denotes the subtree containment relation and  $\text{vars}(e)$  is the set of expression variables that occur free in  $e$ . With an abuse of notation we will use  $\text{vars}(p)$  to denote the capture variables occurring in  $p$  (thus,  $\text{vars}(fe) = \text{vars}(f) \cup \text{vars}(e)$  and  $\text{vars}(p \Rightarrow f) = \text{vars}(f) \setminus \text{vars}(p)$ , the rest of the definition being standard).

As an aside notice that the fact that for a given filter  $f$  the type inference algorithm terminates on all possible input types *does not* imply that the execution of  $f$  terminates on all possible input values. For instance, our analysis correctly detects that for the filter  $\mu X.(x_1, x_2) \Rightarrow X(x_1, x_2) \mid \_ \Rightarrow \_$  type inference terminates on all possible input types (by returning the very same input type) although the application of this same filter never terminates on arguments of the form  $(v_1, v_2)$ .

We can now formally define the two phases of our analysis algorithm.

**First phase.** The first phase is implemented by the function  $\text{Trees}_X(\_, \_)$ . For every filter recursion variable  $X$ , this function explores a filter and does the following two things:

1. It builds a substitution  $\sigma : \mathbf{Vars} \rightarrow \mathbf{Ids}$  from expression variables to variables identifiers, thus associating each capture variable occurring in a pattern of the filter to a fresh identifier for the abstract subtree of the input type it will be bound to.
2. It uses the substitution  $\sigma$  to compute the set of symbolic arguments of recursive calls of  $X$ .

In other words, if  $n$  recursive calls of  $X$  occur in the filter  $f$ , then  $\text{Trees}_X(\sigma, f)$  returns the set  $\{A_1, \dots, A_n\}$  of the  $n$  symbolic arguments of these calls, obtained under the hypothesis that the free variables of  $f$  are associated to subtrees as specified by  $\sigma$ . The formal definition is as follows:

$$\begin{aligned} \text{Trees}_X(\sigma, e) &= \emptyset \\ \text{Trees}_X(\sigma, p \Rightarrow f) &= \text{Trees}_X(\sigma \cup \bigcup_{x_i \in \text{vars}(p)} \{x_i \mapsto \iota_i\}, f) \quad (\iota_i \text{ fresh}) \\ \text{Trees}_X(\sigma, (f_1, f_2)) &= \text{Trees}_X(\sigma, f_1) \cup \text{Trees}_X(\sigma, f_2) \\ \text{Trees}_X(\sigma, f_1 \mid f_2) &= \text{Trees}_X(\sigma, f_1) \cup \text{Trees}_X(\sigma, f_2) \\ \text{Trees}_X(\sigma, \mu X.f) &= \emptyset \\ \text{Trees}_X(\sigma, \mu Y.f) &= \text{Trees}_X(\sigma, f) \quad (X \neq Y) \\ \text{Trees}_X(\sigma, Xa) &= \{a\sigma\} \\ \text{Trees}_X(\sigma, f_1; f_2) &= \text{Trees}_X(\sigma, f_2) \\ \text{Trees}_X(\sigma, of) &= \text{Trees}_X(\sigma, f) \quad (o = \text{groupby}, \text{orderby}) \\ \text{Trees}_X(\sigma, \{\ell_i: f_i, \dots\}_{i \in I}) &= \bigcup_{i \in I} \text{Trees}_X(\sigma, f_i) \end{aligned}$$

where  $a\sigma$  denotes the application of the substitution  $\sigma$  to  $a$ .

The definition above is mostly straightforward. The two important cases are those for the pattern filter where the substitution  $\sigma$  is updated by associating every capture variable of the pattern with a fresh variable identifier, and the one for the recursive call where the symbolic argument is computed by applying  $\sigma$  to the actual argument.

**Second phase.** The second phase is implemented by the function  $\text{Check}()$ . The intuition is that  $\text{Check}(f)$  must “compute” the application of  $f$  to all the symbolic arguments collected in the first phase and then check whether the variables occurring in the arguments of recursive calls (*ie*, the “marked” variables) are actually bound to subtrees of the original type (*ie*, they are bound either to variable identifiers or to subparts of variable identifiers). If we did not have filter composition, then this would be a relatively easy task: it would amount to compute substitutions (by matching symbolic arguments against patterns), apply them, and finally verify that the marked variables satisfy the sought property. Unfortunately, in the case of a composition filter  $f_1; f_2$  the analysis is more complicated than that. Imagine that we want to check the property for the application of  $f_1; f_2$  to a symbolic argument  $A$ . Then to check the property for the recursive calls occurring in  $f_2$  we must compute (or at least, approximate) the set of the symbolic arguments that will be produced by the application of  $f_1$  to  $A$  and that will be thus fed into  $f_2$  to compute the composition. Therefore  $\text{Check}(f)$  will be a function that rather than returning just true or false, it will return a set of symbolic arguments that are the result of the execution of  $f$ , or it will fail if any recursive call does not satisfy the property for marked variables.

More precisely, the function  $\text{Check}()$  will have the form

$$\text{Check}_V(\sigma, f, \{A_1, \dots, A_n\})$$

where  $V \subseteq \mathbf{Vars}$  stores the set of marked variables,  $f$  is a filter,  $\sigma$  is a substitution for (at least) the free variables of  $f$  into  $\mathbf{Ids}$ , and  $A_i$  are symbolic arguments. It will either fail if the marked variables do not satisfy the property of being bound to (subcomponents of) variable identifiers or return an over-approximation of the result of applying  $f$  to all the  $A_i$  under the hypothesis  $\sigma$ . This approximation is either a set of new symbolic arguments, or  $\perp$ . The latter simply indicates that  $\text{Check}()$  is not able to compute the result of the application, typically because it is the result of some expression

- $Check_V(\sigma, f, \{A_1, \dots, A_n\}) = \bigcup_{i=1}^n Check_V(\sigma, f, A_i)$

- **If**  $(\bigcap A[\& \bigcap f] = \emptyset)$  **then**  $Check_V(\sigma, f, A) = \emptyset$  **otherwise:**

$$\begin{aligned}
Check_V(\sigma, e, A) &= \begin{cases} \{a\sigma\} & \text{if } e \equiv a \\ \perp & \text{otherwise} \end{cases} \\
Check_V(\sigma, f_1; f_2, A) &= Check_V(\sigma, f_2, Check_V(\sigma, f_1, A)) \\
Check_V(\sigma, Xa, A) &= \perp \\
Check_V(\sigma, \mu X.f, A) &= \begin{cases} fail & \text{if } Check_{Mark_X}(f)(\sigma, f, A) = fail \\ \perp & \text{otherwise} \end{cases} \\
Check_V(\sigma, f_1|f_2, A) &= \begin{cases} Check_V(\sigma, f_2, A) & \text{if } \bigcap f_1 \bigcap \& \bigcap A = \emptyset \\ Check_V(\sigma, f_1, A) & \text{if } \bigcap A \setminus \bigcap f_1 = \emptyset \\ Check_V(\sigma, f_1, A) \cup Check_V(\sigma, f_2, A) & \text{otherwise} \end{cases} \\
Check_V(\sigma, (f_1, f_2), A) &= \begin{cases} Check_V(\sigma, f_1, A_1) \times Check_V(\sigma, f_2, A_2) & \text{if } A \equiv (A_1, A_2) \\ Check_V(\sigma, f_1, \iota_1) \times Check_V(\sigma, f_2, \iota_2) & \text{if } A \equiv \iota \quad (\text{where } \iota_1 \text{ and } \iota_2 \text{ are fresh}) \\ fail & \text{if } A \equiv \perp \end{cases} \\
Check_V(\sigma, \{\ell_i:f_i, ..\}_{1 \leq i \leq n}, A) &= \begin{cases} \bigcup_{B_i \in Check_V(\sigma, f_i, A_i)} \{\ell_1:B_1, \dots, \ell_n:B_n, \dots, \ell_{n+k}:A_{n+k}, ..\} & \text{if } A \equiv \{\ell_i:f_i, ..\}_{1 \leq i \leq n+k} \\ \bigcup_{B_i \in Check_V(\sigma, f_i, \iota_i)} \{\ell_1:B_1, \dots, \ell_n:B_n, \dots, \ell_{n+k}:A_{n+k}, ..\} & \text{if } A \equiv \iota \quad (\text{where all } \iota_i \text{ are fresh}) \\ fail & \text{if } A \equiv \perp \end{cases} \\
Check_V(\sigma, p \Rightarrow f, A) &= \begin{cases} fail & \text{if } A \not\models p = fail \\ Check_V(\sigma \cup (A \not\models p), f, A) & \text{otherwise} \end{cases}
\end{aligned}$$

**Figure 3.** Definition of the second phase of the analysis

belonging to the host language (eg, the application of a function) or because it is the result of a recursive filter or of a recursive call of some filter.

The full definition of  $Check\_(\_, \_, \_)$  is given in Figure 3. Let us comment the different cases in detail.

$$Check_V(\sigma, f, \{A_1, \dots, A_n\}) = \bigcup_{i=1}^n Check_V(\sigma, f, A_i)$$

simply states that to compute the filter  $f$  on a set of symbolic arguments we have to compute  $f$  on each argument and return the union of the results. Of course, if any  $Check_V(\sigma, f, A_i)$  fails, so does  $Check_V(\sigma, f, \{A_1, \dots, A_n\})$ . The next case is straightforward: if we know that an argument of a given form makes the filter  $f$  fail, then we do not perform any check (the recursive call of  $f$  will never be called) and no result will be returned. For instance, if we apply the filter  $\{\ell:f, ..\}$  to the symbolic argument  $(\iota_1, \iota_2)$ , then this will always fail so it is useless to continue checking this case. Formally, what we do is to consider  $\bigcap A$ , the set of all values that have the form of  $A$  (this is quite simply obtained from  $A$  by replacing any for every occurrence of a variable identifier or of  $\perp$ ) and check whether in it there is any value accepted by  $f$ , namely, whether the intersection  $\bigcap A[\& \bigcap f]$  is empty. If it is so, we directly return the empty set of symbolic arguments, otherwise we have to perform a fine grained analysis according to the form of the filter.

If the filter is an expression, then there are two possible cases: either the expression has the form of an argument ( $\equiv$  denotes

syntactic equivalence), in which case we return it after having applied the substitution  $\sigma$  to it; or it is some other expression, in which case the function says that it is not able to compute a result and returns  $\perp$ .

When the filter is a composition of two filters,  $f_1; f_2$ , we first call  $Check_V(\sigma, f_1, A)$  to analyze  $f_1$  and compute the result of applying  $f_1$  to  $A$ , and then we feed this result to the analysis of  $f_2$ :

$$Check_V(\sigma, f_1; f_2, A) = Check_V(\sigma, f_2, Check_V(\sigma, f_1, A))$$

When the filter is recursive or a recursive call, then we are not able to compute the symbolic result. However, in the case of a recursive filter we must check whether the type inference terminates on it. So we mark the variables of its recursive calls, check whether its definition passes our static analysis, and return  $\perp$  only if this analysis —ie,  $Check_{Mark_X}(f)(\sigma, f, A)$ — did not fail. Notice that  $\perp$  is quite different from failure since it allows us to compute the approximation as long as the filter after the composition does not bind (parts of) the result that are  $\perp$  to variables that occur in arguments of recursive calls. A key example is the case of the filter  $Filter[F]$  defined in Section 5.1. If the parameter filter  $F$  is recursive, then without  $\perp$   $Filter[F]$  would be rejected by the static analysis. Precisely, this is due to the composition  $Fx;('false \Rightarrow X| \dots)$ . When  $F$  is recursive the analysis supposes that  $Fx$  produces  $\perp$ , which is then passed over to the pattern. The recursive call  $Xl$  is executed in the environment  $\perp / 'true$ , which is empty. Therefore

$$\begin{aligned}
A/x &= \begin{cases} \{x \mapsto A\} & \text{if } A \equiv v \text{ or } A \equiv \iota \text{ or } x \notin \mathcal{V} \\ \text{fail} & \text{otherwise} \end{cases} \\
A/(p_1, p_2) &= \begin{cases} (A_1/p_1) \cup (A_2/p_2) & \text{if } A \equiv (A_1, A_2) \\ (\iota_1/p_1) \cup (\iota_2/p_2) & \text{if } A \equiv \iota \quad (\text{where } \iota_1, \iota_2 \text{ are fresh}) \\ \text{fail} & \text{if } A \equiv \perp \\ \Omega & \text{otherwise} \end{cases} \\
A/p_1 \&p_2 &= A/p_1 \cup A/p_2 \\
A/p_1|p_2 &= \begin{cases} A/p_1 & \text{if } A/p_1 \neq \Omega \\ A/p_2 & \text{otherwise} \end{cases} \\
A/\{\ell_1:p_1, \dots, \ell_n:p_n\} &= \begin{cases} \bigcup_{i \leq n} A_i/p_i & \text{if } A \equiv \{\ell_1:A_1, \dots, \ell_n:A_n\} \\ \bigcup_{i \leq n} \iota_i/p_i & \text{if } A \equiv \iota \quad (\text{where } \iota_i \text{ are fresh}) \\ \text{fail} & \text{if } A \equiv \perp \\ \Omega & \text{otherwise} \end{cases} \\
A/\{\ell_1:p_1, \dots, \ell_n:p_n, ..\} &= \begin{cases} \bigcup_{i \leq n} A_i/p_i & \text{if } A \equiv \{\ell_1:A_1, \dots, \ell_n:A_n, \dots, \ell_{n+k}:A_{n+k}\} \\ \bigcup_{i \leq n} \iota_i/p_i & \text{if } A \equiv \iota \quad (\text{where } \iota_i \text{ are fresh}) \\ \text{fail} & \text{if } A \equiv \perp \\ \Omega & \text{otherwise} \end{cases} \\
A/t &= \emptyset
\end{aligned}$$

**Figure 4.** Pattern matching with respect to a set  $\mathcal{V}$  of marked variables

the result of the  $Fx$  call, whatever it is, cannot affect the termination of the subsequent recursive calls. Even though the composition uses the result of  $Fx$ , the form of this result cannot be such as to make typechecking of  $\text{Filter}[F]$  diverge.

For the union filter  $f_1|f_2$  there are three possible cases:  $f_1$  will always fail on  $A$ , so we return the result of  $f_2$ ; or  $f_2$  will never be executed on  $A$  since  $f_1$  will never fail on  $A$ , so we return just the result of  $f_1$ ; or we cannot tell which one of  $f_1$  or  $f_2$  will be executed, and so we return the union of the two results.

If the filter is a product filter  $(f_1, f_2)$ , then its accepted type  $\mathcal{I}((f_1, f_2))$  is a subtype of  $(\text{any}, \text{any})$ . Since we are in the case where  $\mathcal{I}A[\&\mathcal{I}((f_1, f_2))]$  is not empty, then  $A$  can have just two forms: either it is a pair or it is a variable identifier  $\iota$ . In the former case we check the filters component-wise and return as result the set-theoretic product of the results. In the latter case, recall that  $\iota$  represents a subtree of the original input type. So if the application does not fail it means that each subfilter will be applied to a subtree of  $\iota$ . We introduce two fresh variable identifiers to denote these subtrees of  $\iota$  (which, of course are subtrees of the original input type, too) and, as in the previous case, apply the check component-wise and return the set-theoretic product of the results. The case for the record filter is similar to the case of the product filter.

Finally, the case of pattern filter is where the algorithm checks that the marked variables are indeed bound to subtrees of the initial input type. What  $\text{Check}()$  does is to match the symbolic argument against the pattern and update the substitution  $\sigma$  so as to check the subfilter  $f$  in an environment with the corresponding assignments for the capture variables in  $p$ . Notice however that the pattern matching  $A \vee p$  receives as extra argument the set  $\mathcal{V}$  of marked variables, so that while computing the substitutions for the capture

variables in  $p$  it also checks that all capture variables that are also marked are actually bound to a subtree of the initial input type.  $A \vee p$  is defined in Figure 4 (to enhance readability we omitted the  $\mathcal{V}$  index since it matters only in the first case and it does not change along the definition). The check for marked variables is performed in the first case of the definition: when a symbolic argument  $A$  is matched against a variable  $x$ , if the variable is not marked ( $x \notin \mathcal{V}$ ), then the substitution  $\{x \mapsto A\}$  is returned; if it is marked, then the matching (and thus the analysis) does not fail only if the symbolic argument is either a value (so it does not depend on the input type) or it is a variable identifier (so it will be exactly a subtree of the original input type). The other cases of the definition of  $A \vee p$  are standard. Just notice that in the product and record patterns when  $A$  is a variable identifier  $\iota$ , the algorithm creates fresh new variable identifiers to denote the new subtrees in which  $\iota$  will be decomposed by the pattern. Finally, the reader should not confound  $\Omega$  and  $\text{fail}$ . The former indicates that the argument does not match the value, while the latter indicates that a marked variable is not bound to a value or to exactly a subtree of the initial input type (notice that the case  $\Omega$  cannot happen in the definition of  $\text{Check}()$  since  $A \vee p$  is called only when  $\mathcal{I}A$  is contained in  $\mathcal{I}p$ ).

**The two phases together.** Finally we put the two phases together. Given a recursive filter  $\mu X.f$  we run the analysis by marking the variables in the recursive calls of  $X$  in  $f$ , running the first phase and then feeding the result to the second phase:  $\text{Check}_{\text{Mark}_X(f)}(\emptyset, f, \text{Trees}_X(\emptyset, f))$ . If this function does not fail then we say that  $\mu X.f$  passed the check:

$$\text{Check}(\mu X.f) \iff \text{Check}_{\text{Mark}_X(f)}(\emptyset, f, \text{Trees}_X(\emptyset, f)) \neq \text{fail}$$



For filters that are not recursive at toplevel it suffices to add a dummy recursion:  $Check(f) \stackrel{\text{def}}{=} Check(\mu X.f)$  with  $X$  fresh.

**Theorem 16 (Termination).** *If  $Check(f)$  then the type inference algorithm terminates for  $f$  on every input type  $t$ .*

In order to prove the theorem above we need some auxiliary definitions:

**Definition 17 (Plinth [15]).** A plinth  $\sqsupset \subset \mathbf{Types}$  is a set of types with the following properties:

- $\sqsupset$  is finite
- $\sqsupset$  contains any, empty and is closed under Boolean connectives ( $\&, |, \neg$ )
- for all types  $t = (t_1, t_2)$  in  $\sqsupset$ ,  $t_1 \in \sqsupset$  and  $t_2 \in \sqsupset$
- for all types  $t = \{\ell_1:t_1, \dots, \ell_n:t_n, (..)\}$  in  $\sqsupset$ ,  $t_i \in \sqsupset$  for all  $i \in [1..n]$ .

We define the plinth of  $t$ , noted  $\sqsupset(t)$ , as the smallest plinth containing  $t$ .

Intuitively, the plinth of  $t$  is the set of types that can be obtained by all possible boolean combination of the subtrees of  $t$ . Notice that  $\sqsupset(t)$  is always defined since our types are regular: they have finitely many distinct subtrees, which (modulo type equivalence) can thus be combined in finitely many different ways.

**Definition 18 (Extended plinth and support).** Let  $t$  be a type and  $f$  a filter. The *extended plinth* of  $t$  and  $f$ , noted  $\hat{\sqsupset}(f, t)$  is defined as  $\hat{\sqsupset}(t \vee \bigvee_{v \sqsubseteq f} v)$  (where  $v$  ranges over values).

The *support* of  $t$  and  $f$ , noted as  $Support(f, t)$ , is defined as

$$Support(f, t) \stackrel{\text{def}}{=} \hat{\sqsupset}(f, t) \cup \bigcup_{A \in Trees_{\emptyset}(\emptyset, f)} \{A\sigma \mid \sigma : \mathbf{Ids}(A) \rightarrow \hat{\sqsupset}(f, t)\}$$

The extended plinth of  $t$  and  $f$  is the set of types that can be obtained by all possible boolean combination of the subtrees of  $t$  and of values that occur in the filter  $f$ . The intuition underlying the definition of support of  $t$  and  $f$  is that it includes all the possible types of arguments of recursive calls occurring in  $f$ , when  $f$  is applied to an argument of type  $t$ .

Lastly, let us prove the following technical lemma:

**Lemma 19.** *Let  $f$  be a filter such that  $Check(f)$  holds. Let  $t$  be a type. For every derivation  $D$  (finite or infinite) of*

$$\Gamma'; \Delta'; M' \vdash_{\text{fil}} f(t) : s$$

*for every occurrence of the rule [FIL-CALL-NEW]*

$$\frac{\Gamma'; \Delta'; M', ((X, t'') \mapsto T) \vdash_{\text{fil}} \Delta'(X)(t'') : t' \quad \begin{matrix} t'' = \text{type}(\Gamma', a) \\ (X, t'') \notin \text{dom}(M') \\ T \text{ fresh} \end{matrix}}{\Gamma'; \Delta'; M' \vdash_{\text{fil}} (X a)(s) : \mu T.t'}$$

*in  $D$ , for every  $x \in \text{vars}(a)$ ,  $\Gamma'(x) \in \hat{\sqsupset}(f, t)$  (or equivalently,  $\text{type}(\Gamma', a) \in Support(f, t)$ )*

*Proof.* By contradiction. Suppose that there exists an instance of the rule for which  $x \in \text{vars}(a) \wedge x \notin \hat{\sqsupset}(f, t)$ . This means that  $\Gamma'(x)$  is neither a singleton type occurring in  $f$  nor a type obtained from  $t$  by applying left projection, right projection or label selection ( $\star$ ). But since  $Check(f)$  holds,  $x$  must be bound to either an identifier or a value  $v$  (see Figure 4) during the computation of  $Check_{\text{Mark}_X(f)}(\emptyset, f, Trees_X(\emptyset, f))$ . Since identifiers in  $Check(f)$  are only introduced for the input parameter or when performing a left projection, right projection or label selection of another identifier, this ensure that  $x$  is never bound to the result of an expression whose type is not in  $\hat{\sqsupset}(f, t)$ , which contradicts ( $\star$ ).  $\square$

We are now able to prove Theorem 16. Let  $t$  be a type and  $f$  a filter define

$$\mathcal{D}(f, t) = \{(X, t') \mid X \sqsubseteq f, t' \in Support(f, t)\}$$

Notice that since  $Support(f, t)$  is finite and there are finitely many different recursion variables occurring in a filter, then  $\mathcal{D}(f, t)$  is finite, too. Now let  $f$  and  $t$  be the filter and type mentioned in the statement of Theorem 16 and consider the (possibly infinite) derivation of  $\Gamma'; \Delta'; M' \vdash_{\text{fil}} f(t) : s$ . Assign to every judgment  $\Gamma'; \Delta'; M' \vdash_{\text{fil}} f'(t') : s'$  the following weight

$$Wgt(\Gamma'; \Delta'; M' \vdash_{\text{fil}} f'(t') : s') = (|\mathcal{D}(f, t) \setminus \text{dom}(M')|, \text{size}(f'))$$

where  $|S|$  denotes the cardinality of the set  $S$  (notice that in the definition  $S$  is finite),  $\text{dom}(M')$  is the domain of  $M'$ , that is, the set of pairs  $(X, t)$  for which  $M'$  is defined, and  $\text{size}(f')$  is the depth of the syntax tree of  $f'$ .

Notice that the set of all weights lexicographically ordered form a well-founded order. It is then not difficult to prove that every application of a rule in the derivation of  $\Gamma'; \Delta'; M' \vdash_{\text{fil}} f(t) : s$  strictly decreases  $Wgt$ , and therefore that this derivation must be finite. This can be proved by case analysis on the applied rule, and we must distinguish only three cases:

[FIL-FIX] Notice that in this case the first component of  $Wgt$  either decreases or remains constant, and the second component strictly decreases.

[FIL-CALL-NEW] In this case Lemma 19 ensures that the first component of the  $Wgt$  of the premise strictly decreases. Since this is the core of the proof let us expand the rule. Somewhere in the derivation of  $\Gamma'; \Delta'; M' \vdash_{\text{fil}} f(t) : s$  we have the following rule:

$$\frac{\text{[FIL-CALL-NEW]} \quad \Gamma'; \Delta'; M', ((X, t'') \mapsto T) \vdash_{\text{fil}} \Delta'(X)(t'') : t' \quad \begin{matrix} t'' = \text{type}(\Gamma', a) \\ (X, t'') \notin \text{dom}(M') \\ T \text{ fresh} \end{matrix}}{\Gamma'; \Delta'; M' \vdash_{\text{fil}} (X a)(s) : \mu T.t'}$$

and we want to prove that  $(\mathcal{D}(f, t) \setminus \text{dom}(M')) \supsetneq (\mathcal{D}(f, t) \setminus (\text{dom}(M') \cup (X, t'')))$  and the containment must be strict to ensure that the measure decreases. First of all notice that  $(X, t'') \notin \text{dom}(M')$ , since it is a side condition for the application of the rule. So in order to prove that containment is strict it suffices to prove that  $(X, t'') \in \mathcal{D}(f, t)$ . But this is a consequence of Lemma 19 which ensures that  $t'' \in Support(f, t)$ , whence the result.

[FIL- $\star$ ] With all other rules the first component of  $Wgt$  remains constant, and the second component strictly decreases.

## A.1 Improvements

Although the analysis performed by our algorithm is already fine grained, it can be further refined in two simple ways. As explained in Section 4, the algorithm checks that after *one step* of reduction the capture variables occurring in recursive calls are bound to subtrees of the initial input type. So a possible improvement consists to try to check this property for a higher number of steps, too. For instance consider the filter:

$$\begin{aligned} \mu X. (& \text{nil}, x) \Rightarrow X(\{\ell : x\}) \\ & | (y, \_) \Rightarrow X((\text{nil}, (y, y))) \\ & | \_ \Rightarrow \text{nil} \end{aligned}$$

This filter does not pass our analysis since if  $\iota_y$  is the identifier bound to the capture variable  $y$ , then when unfolding the recursive call in the second branch the  $x$  in the first recursive call will be bound to  $(\iota_y, \iota_y)$ . But if we had pursued our abstract execution one step further we would have seen that  $(\iota_y, \iota_y)$  is not used in a recursive call and, thus, that type inference terminates. Therefore, a first improvements is to modify  $Check()$  so that it does not stop

```

type M = (K,V) | (V,V,K)
type V = { var : string } | { lambda2 : (string, string, M) }
type K = { var : string } | { lambda1 : (string, M) }

let filter Eval =
  | ( { lambda2 : (x, k, m) }, v , h ) -> m ; Subst[x,v] ; Subst[k,h] ; Eval
  | ( { lambda1 : (x, m) }, v ) -> m ; Subst[x,v] ; Eval
  | x -> x

let filter Subst[c,F] =
  | ( Subst[c,F] , Subst[c,F] , Subst[c,F] )
  | ( Subst[c,F] , Subst[c,F] )
  | { var : c } -> F
  | { lambda2 : (x&-c, k&-c, m) } -> { lambda2 : (x, k, m;Subst[c,F]) }
  | { lambda1 : (x&-c, m) } -> { lambda1 : (x, m;Subst[c,F]) }
  | x -> x

```

Figure 5. Filter encoding of  $\lambda_{cps}$

just after one step of reduction but tries to go down  $n$  steps, with  $n$  determined by some heuristics based on sizes of the filter and of the input type.

A second and, in our opinion, more promising improvement is to enhance the precision of the test  $\lambda A[\& \lambda f] = \emptyset$  in the definition of *Check*, that verifies whether the filter  $f$  fails on the given symbolic argument  $A$ . In the current definition the only information we collect about the type of symbolic arguments is their structure. But further type information, currently unexploited, is provided by patterns. For instance, in the following (admittedly stupid) filter

$$\begin{aligned} \mu X.(\text{int}, x) \Rightarrow x \\ | y \& \text{int} \Rightarrow X((y, y)) \\ | z \Rightarrow z \end{aligned}$$

if in the first pass we associate  $y$  with  $\iota_y$ , then we know that  $(\iota_y, \iota_y)$  has type  $(\text{int}, \text{int})$ . If we record this information, then we know that in the second pass  $(\iota_y, \iota_y)$  will always match the first pattern, and so it will never be the argument of a further recursive call. In other words, there is at most one nested recursive call. The solution is conceptually simple (but yields a cumbersome formalization, which is why we chose the current formulation) and amounts to modify *Trees*( $\cdot$ ) so that when it introduces fresh variables it records their type information with them. It then just suffices to modify the definition of  $\lambda A[\& \lambda f]$  so as it is obtained from  $A$  by replacing every occurrence of a variable identifier by its type information (rather than by *any*) and the current definition of *Check*() will do the rest.

## B. Proof of Turing completeness (Theorem 7)

In order to prove Turing completeness we show how to define by filters an evaluator for untyped (call-by-value)  $\lambda$ -calculus. If we allowed recursive calls to occur on the left-hand side of composition, then the encoding would be straightforward: just implement  $\beta$  and context reductions. The goal however is to show that the restriction on compositions does not affect expressiveness. To that end we have to avoid context reductions, since they require recursive calls before composition. To do so, we first translate  $\lambda$ -terms via Plotkin's call-by-value CPS translation and apply Steele and Rabbit's administrative reductions to them obtaining terms in  $\lambda_{cps}$ . The latter is isomorphic to cbv  $\lambda$ -calculus (see for instance [30]) and defined as follows.

$$\begin{aligned} M &::= KV \mid VVK \\ V &::= x \mid \lambda x. \lambda k. M \\ K &::= k \mid \lambda x. M \end{aligned}$$

with the following reduction rules (performed at top-level, without any reduction under context).

$$\begin{aligned} (\lambda x. \lambda k. M)VK &\longrightarrow M[x := V][k := K] \\ (\lambda x. M)V &\longrightarrow M[x := V] \end{aligned}$$

In order to prove Turing completeness it suffices to show how to encode  $\lambda_{cps}$  terms and reductions in our calculus. For the sake of readability, we use mutually recursive types (rather than their encoding in  $\mu$ -types), we use records (though pairs would have sufficed), we write just the recursion variable  $X$  for the filter  $x \rightarrow Xx$ , and use  $\neg t$  to denote the type *any*  $\setminus t$ . Term productions are encoded by the recursive types given at the beginning of Figure 5.

Next we define the evaluation filter *Eval*. In its body it calls the filter *Subst*[ $c, F$ ] which implements the capture free substitution and, when  $c$  denotes a constant, is defined as right below.

It is straightforward to see that the definitions in Figure 5 implement the reduction semantics of CPS terms. Of course the definition above would not pass our termination condition. Indeed while *Subst* would be accepted by our algorithm, *Eval* would fail since it is not possible to ensure that the recursive calls of *Eval* will receive from *Subst* subtrees of the original input type. This is expected: while substitution always terminate they may return trees that are not subtrees of the original term.

## C. Proof of subject reduction (Theorem 13)

We first give the proof of Lemma 9, which we restate:

Let  $f$  be a filter and  $v$  be a value such that  $v \notin \lambda f$ . Then for every  $\gamma, \delta$ , if  $\delta; \gamma \vdash_{eval} f(v) \rightsquigarrow r$ ,  $r \equiv \Omega$ .

*Proof.* By induction on the derivation of  $\delta; \gamma \vdash_{eval} f(v) \rightsquigarrow r$ , and case analysis on the last rule of the derivation:

- (**expr**): here,  $\lambda e$  = *any* for any expression  $e$ , therefore this rule cannot be applied (since  $\forall v, v \in \text{any}$ ).
- (**prod**): assume  $v \notin \lambda(f_1, f_2) \equiv (\lambda f_1, \lambda f_2)$ . then either  $v \notin (\text{any}, \text{any})$ , in which case only rule (**error**) applies and therefore  $r \equiv \Omega$ . Or  $v \equiv (v_1, v_2)$  and  $v_1 \notin \lambda f_1$ . Then by induction hypothesis on the first premise,  $\delta; \gamma \vdash_{eval} f_1(v_1) \rightsquigarrow r_1$  and  $r_1 \equiv \Omega$ , which contradicts the side condition  $r_1 \neq \Omega$  (and

We were a little bit sloppy in the notation and used a filter parameter as a pattern. Strictly speaking this is not allowed in filters and, for instance, the branch  $\{ \text{var} : c \} \rightarrow F$  in *Subst* should rather be written as  $\{ \text{var} : y \} \rightarrow ( (y=c); (\text{true} \rightarrow F \mid \text{false} \rightarrow \{ \text{var} : y \}) )$ . Similarly, for the other two cases.

similarly for the second premise). Therefore this rule cannot be applied to evaluate  $(f_1, f_2)$ .

**(patt)**: Similarly to the previous case. If  $v \notin \lambda p \Rightarrow f \hat{=} \lambda p \hat{\&} \lambda f \hat{}$  then either  $v \notin \lambda p \hat{}$  (which contradicts  $v/p \neq \Omega$ ) or  $v \notin \lambda f \hat{}$  and by induction hypothesis,  $\delta; \gamma, v/p \vdash_{eval} f(v) \rightsquigarrow \Omega$ .

**(comp)**: If  $v \notin \lambda f_1; f_2 \hat{=} \lambda f_1 \hat{}$ , then by induction hypothesis on the first premise,  $r_1 \equiv \Omega$ , which contradicts the side condition  $r_1 \neq \Omega$ . Therefore only rule **(error)** can be applied here.

**(recd)**: Similar to product type: either  $v$  is not a record and therefore only the rule **(error)** can be applied, or one of the  $r_i = \Omega$  (by induction hypothesis) which contradicts the side condition.

**(union1)** and **(union2)**: since  $v \notin \lambda f_1 | f_2 \hat{}$ , this means that  $v \notin \lambda f_1 \hat{}$  and  $v \notin \lambda f_2 \hat{}$ . Therefore if rule **(union1)** is chosen, by induction hypothesis,  $r_1 \equiv \Omega$  which contradicts the side condition. If rule **(union2)** is chosen, then by induction hypothesis  $r_2 \equiv \Omega$  which gives  $r \equiv \Omega$ .

**(rec-call)** trivially true, since it cannot be that  $v \notin \text{any}$

**(rec)** we can apply straightforwardly the induction hypothesis on the premise and have that  $r \equiv \Omega$ .

**(groupby)** and **(orderby)**: If  $v \notin f$  then the only rule that applies is **(error)** which gives  $r \equiv \Omega$ .

□

We are now equipped to prove the subject reduction theorem which we restate in a more general manner:

For every  $\Gamma, \Delta, M, \gamma, \delta$  such that  $\forall x \in \text{dom}(\gamma), x \in \text{dom}(\Gamma) \wedge \gamma(x) : \Gamma(x)$ , if  $\Gamma; \Delta; M \vdash_{\text{fl}} f(t) : s$ , then for all  $v : t$ ,  $\delta; \gamma \vdash_{eval} f(v) \rightsquigarrow r$  implies  $r : s$ .

The proof is by induction on the derivation of  $\delta; \gamma \vdash_{eval} f(v) \rightsquigarrow r$ , and by case analysis on the rule. Beside the basic case, the only rules which are not a straightforward application of induction are the rules **union1** and **union2** that must be proved simultaneously. Other cases being proved by a direct application of the induction hypothesis, we only detail the case of the product constructor.

**(expr)**: we suppose that the host languages enjoys subject reduction, hence  $e(v) = r : s$ .

**(prod)**: Here, we know that  $t \equiv \bigvee_{i \leq \text{rank}(t)} (t_1^i, t_2^i)$ . Since  $v$  is a value, and  $v : t$ , we obtain that  $v \equiv (v_1, v_2)$ . Since  $(v_1, v_2)$  is a value,  $\exists i \leq \text{rank}(t)$  such that  $(v_1, v_2) : (t_1^i, t_2^i)$ . The observation that  $v$  is a value is crucial here, since in general given a type  $t' \leq t$  with  $t' \equiv \bigvee_{j \leq \text{rank}(t')} (t_1'^j, t_2'^j)$  it is not true that  $\exists i, j$  s.t.  $(t_1'^j, t_2'^j) \leq (t_1^i, t_2^i)$ . However this property holds for singleton types and therefore for values. We have therefore that  $v_1 : t_1^i$  and  $v_2 : t_2^i$ . Furthermore, since we suppose that a typing derivation exists and the typing rules are syntax directed, then  $\Gamma; \Delta; M \vdash_{\text{fl}} f(t_1^i) : s_1^i$  must be a used to prove our typing judgement. We can apply the induction hypothesis and deduce that  $\delta; \gamma \vdash_{eval} f(v_1) \rightsquigarrow r_1$  and similarly for  $v_2$ . We have therefore that the filter evaluates to  $(r_1, r_2)$  which has type  $(s_1^i, s_2^i) \leq \bigvee_{j \leq \text{rank}(s)} (s_1^j, s_2^j)$  which proves this case.

**(union1)** and **(union2)**: both rules must be proved together. Indeed, given a filter  $f_1 | f_2$  for which we have a typing derivation for the judgement  $\Gamma; \Delta; M \vdash_{\text{fl}} f_1 | f_2(t) : s$ , either  $v : t \hat{\&} \lambda f_1 \hat{}$  and we can apply the induction hypothesis and therefore,  $\Gamma; \Delta; M \vdash_{\text{fl}} f_1 | (t \hat{\&} \lambda f_1 \hat{}) : s_1$  and if  $\delta; \gamma \vdash_{eval} f(v_1) \rightsquigarrow r_1$  (case **union1**)  $r_1 : s_1$ . However it might be that  $v \notin t \hat{\&} \lambda f_1 \hat{}$ . Then by Lemma 9 we have that  $\delta; \gamma \vdash_{eval} f(v_1) \rightsquigarrow \Omega$  (case **union2**) and we can apply the induction hypothesis on the second premise, which gives us  $\delta; \gamma \vdash_{eval} f(v_1) \rightsquigarrow r_2$  which allows us to conclude.

**(error)**: this rule can never occur. Indeed, if  $v : t$  and  $\Gamma; \Delta; M \vdash_{\text{fl}} f(t) : s$ , that means in particular that  $t \leq \lambda f \hat{}$  and therefore that all of the side conditions for the rules **prod**, **pat** and **recd** hold.

## D. Complexity of the typing algorithm

*Proof.* For clarity we restrict ourselves to types without record constructors but the proof can straightforwardly be extended to them. First remark that our types are isomorphic to alternating tree automata (ATA) with intersection, union and complement (such as those defined in [10, 17]). From such an ATA  $t$  it is possible to compute a non-deterministic tree automaton  $t'$  of size  $O(2^{|t|})$ . When seeing  $t'$  in our type formalism, it is a type generated by the following grammar:

$$\begin{aligned} \tau &::= \mu X. \tau \mid \tau | \text{const} \mid \text{const} \\ \text{const} &::= (\tau, \tau) \mid \text{atom} \\ \text{atom} &::= X \mid b \end{aligned}$$

where  $b$  ranges over negation and intersections of basic types. Intuitively each recursion variable  $X$  is a state in the NTA, a finite union of products is a set of non-deterministic transitions whose right-hand-side are each of the products and *atom* productions are leaf-states. Then it is clear that if  $\text{Check}(f)$  holds, the algorithm considers at most  $|f| \times |t'|$  distinct cases (thanks to the memoization set  $M$ ). Furthermore for each rule, we may test a subtyping problem (eg,  $t \leq \lambda f \hat{}$ ), which itself is EXPTIME thus giving the bound □

## E. Precise typing of groupby

The process of inferring a precise type for **groupby**  $f(t)$  is decomposed in several steps. First we have to compute the set  $\mathcal{D}(f)$  of discriminant domains for the filter  $f$ . The idea is that these domains are the types on which we know that  $f$  may give results of different types. Typically this corresponds to all possible branchings that the filter  $f$  can do. For instance, if  $\{t_1, \dots, t_n\}$  are pairwise disjoint types and  $f$  is the filter  $t_1 \Rightarrow e_1 \mid \dots \mid t_n \Rightarrow e_n$ , then the set of discriminant domains for the filter  $f$  is  $\{t_1, \dots, t_n\}$ , since the various  $s_i$  obtained from  $\emptyset; \emptyset; \emptyset \vdash_{\text{fl}} f(t_i) : s_i$  may all be different, and we want to keep track of the relation between a result type  $s_i$  and the input type  $t_i$  that produced it. Formally  $\mathcal{D}(f)$  is defined as follows:

$$\begin{aligned} \mathcal{D}(e) &= \{\text{any}\} \\ \mathcal{D}((f_1, f_2)) &= \mathcal{D}(f_1) \times \mathcal{D}(f_2) \\ \mathcal{D}(p \Rightarrow f) &= \{\lambda p \hat{\&} t \mid t \in \mathcal{D}(f)\} \\ \mathcal{D}(f_1 | f_2) &= \mathcal{D}(f_1) \cup \{t \mid \lambda f_1 \hat{} \mid t \in \mathcal{D}(f_2)\} \\ \mathcal{D}(\mu X. f) &= \mathcal{D}(f) \\ \mathcal{D}(X a) &= \{\text{any}\} \\ \mathcal{D}(f_1; f_2) &= \mathcal{D}(f_1) \\ \mathcal{D}(o f) &= \{[\text{any}^*]\} \quad (o = \text{groupby}, \text{orderby}) \\ \mathcal{D}(\{ \ell_i : f_i, .. \}_{i \in I}) &= \bigcup_{t_i \in \mathcal{D}(f_i)} \{ \ell_i : t_i, .. \}_{i \in I} \end{aligned}$$

Now in order to have a high precision in typing we want to compute the type of the intermediate list of the **groupby** on a set of disjoint types. For that we define a normal form of a set of types that given a set of types computes a new set formed of pairwise disjoint types whose union covers the union of all the original types.

$$\mathcal{N}(\{t_i \mid i \in S\}) = \bigcup_{\emptyset \subset I \subseteq S} \left( \bigcap_{i \in I} t_i \setminus \bigcup_{j \in S \setminus I} t_j \right)$$

Recall that we are trying to deduce a type for **groupby**  $f(t)$ , that is for the expression **groupby**  $f$  when applied to an argument of type  $t$ . Which types shall we use as input for our filter  $f$  to compute the type of the intermediate result? Surely we want to have all the

types that are in the  $\text{item}(t)$ . This however is not enough since we would not use the information of the discriminant domains for the filter. For instance if the filter gives two different result types for positive and negative numbers and the input is a list of integers, we want to compute two different result types for positive and negatives and not just to compute the result of the filter application on generic integers (indeed  $\text{item}(\text{int}^*) = \{\text{int}\}$ ). So the idea is to add to the set that must be normalized all the types of the set of discriminant types. This however is not precise enough, since these domains may be much larger than the item types of the input list. For this reason we just take the part of the domain types that intersects at least some possible value of the list. In other terms we consider the normal form of the following set

$$\text{item}(t) \cup \{ \sigma_i \wedge \bigvee_{t_i \in \mathcal{D}(f)} t_i \mid \sigma_i \in \text{item}(t) \}$$

The idea is then that if  $\{t_i\}_{i \in I}$  is the normalized set of the set above, then the type of grouby is the type  $[\bigvee_{i \in I} (f(t_i), [t_i+])^*]$  with the optimization that we can replace the  $*$  by a  $+$  if we know that the input list is not empty.

## F. Comparison with top-down tree transducers

We first show that every Top-down tree transducers with regular look-ahead can be encoded by a filter. We use the definition of [14] for top-down tree transducers:

**Definition 20 (Top-down Tree transducer with regular look-ahead).** A top-down tree transducer with regular look-ahead (TDTTR) is a 5-tuple  $\langle \Sigma, \Delta, Q, Q_d, R \rangle$  where  $\Sigma$  is the input alphabet,  $\Delta$  is the output alphabet,  $Q$  is a set of states,  $Q_d$  the set of initial states and  $R$  a set of rules of the form:

$$q(a(x_1, \dots, x_n)), D \rightarrow b(q_1(x_{i_1}), \dots, q_m(x_{i_m}))$$

where  $a \in \Sigma_n$ ,  $b \in \Delta_m$ ,  $q, q_1, \dots, q_m \in Q$ ,  $\forall j \in 1..m, i_j \in 1..n$  and  $D$  is a mapping from  $\{x_1, \dots, x_n\}$  to  $2^{T_\Sigma}$  ( $T_\Sigma$  being the set of regular tree languages over alphabet  $\Sigma$ ). A TDTTR is *deterministic* if  $Q_d$  is a singleton and the left-hand sides of the rules in  $R$  are pairwise disjoint.

Since our filters encode programs, it only make sense to compare filters and *deterministic* TDTTR. We first show that given such a *deterministic* TDTTR we can write an equivalent filter  $f$  and, furthermore, that  $\text{Check}(f)$  does not fail. First we recall the encoding of ranked labeled trees into the filter data-model:

**Definition 21 (Tree encoding).** Let  $t \in T_\Sigma$ . We define the tree-encoding of  $t$  and we write  $\llbracket t \rrbracket$  the value defined inductively by:

$$\begin{aligned} \llbracket a \rrbracket &= ('a, []) & \forall a \in \Sigma_0 \\ \llbracket a(t_1, \dots, t_n) \rrbracket &= ('a, [\llbracket t_1 \rrbracket \dots \llbracket t_n \rrbracket]) & \forall a \in \Sigma_n \end{aligned}$$

where the list notation  $[v_1 \dots v_n]$  is a short-hand for  $(v_1, \dots, (v_n, \text{'nil'}))$ . We generalize this encoding to tree languages and types. Let  $S \subseteq T_\sigma$ , we write  $\llbracket S \rrbracket$  the set of values such that  $\forall t \in S, \llbracket t \rrbracket \in \llbracket S \rrbracket$ .

In particular it is clear than when  $S$  is regular,  $\llbracket S \rrbracket$  is a type.

**Lemma 22 (TDTTR  $\rightarrow$  filters).** Let  $T = \langle \Sigma, \Delta, Q, \{q_0\}, R \rangle$  be a deterministic TDTTR. There exists a filter  $f_T$  such that:

$$\forall t \in \text{dom}(T), \llbracket T(t) \rrbracket \equiv f_T \llbracket t \rrbracket$$

*Proof.* The encoding is as follows. For every state  $q_i \in Q$ , we will introduce a recursion variable  $X_i$ . Formally, the translation is performed by the function  $TR : 2^Q \times Q \rightarrow \mathbf{Filters}$  defined as:

$$\begin{aligned} TR(S, q_i) &= x \Rightarrow X_i x & \text{if } q_i \in S \\ TR(S, q_i) &= \mu X_i. (f_1 \dots f_n) & \text{if } q_i \notin S \end{aligned}$$

where every rule  $r_j \in R$

$$r_j \equiv q_i(a_j(x_1, \dots, x_n)), D_j \rightarrow b_j(q_{j_1}(x_{j_{k_1}}), \dots, q_{j_m}(x_{j_{k_m}}))$$

is translated into:

$$\begin{aligned} ('a_j, [x_1 \& \llbracket D_j(x_1) \rrbracket \dots x_n \& \llbracket D_j(x_n) \rrbracket]) &\Rightarrow 'b_j, [] & \text{if } b_j \in \Delta_0 \\ ('a_j, [x_1 \& \llbracket D_j(x_1) \rrbracket \dots x_n \& \llbracket D_j(x_n) \rrbracket]) &\Rightarrow & \\ ('b_j, (x_{j_{k_1}}; TR(S', q_{j_1}), \dots, (x_{j_{k_m}}; TR(S', q_{j_m}), \text{'nil'}))) && \\ \text{where } S' = S \cup \{q_i\} && \text{otherwise} \end{aligned}$$

The fact that  $\forall t \in \text{dom}(T), \llbracket T(t) \rrbracket \equiv TR(\emptyset, q_0) \llbracket t \rrbracket$  is proved by a straightforward induction on  $t$ . The only important point to remark is that since  $T$  is deterministic, there is exactly one branch of the alternate filter  $f_1 \dots f_n$  that can be selected by pattern matching for a given input  $v$ . For as to why  $\text{Check}(f_T)$  holds, it is sufficient to remark that each recursive call is made on a strict subtree of the input which guarantees that  $\text{Check}(f_T)$  returns true.  $\square$

**Lemma 23 (TDTTR  $\not\leftarrow$  filters).** Filters are strictly more expressive than TDTTRs.

*Proof.* Even if we restrict filters to have the same domain as TDTTRs (meaning, we used a fixed input alphabet  $\Sigma$  for atomic types) we can define a filter that cannot be expressed by a TDTTR. For instance consider the filter:

$$\begin{aligned} y &\Rightarrow (\mu X. ('a, \text{'nil'}) \Rightarrow y \\ &\quad | ('b, \text{'nil'}) \Rightarrow y \\ &\quad | ('a, [x]) \Rightarrow ('a, [X x]) \\ &\quad | ('b, [x]) \Rightarrow ('b, [X x])) \end{aligned}$$

This filter works on monadic trees of the form  $u_1(\dots u_{n-1}(u_n) \dots)$  where  $u_i \in \{a, b\}$ , and essentially replaces the leaf  $u_n$  of a tree  $t$  by a copy of  $t$  itself. This cannot be done by a TDTTR. Indeed, TDTTR have only two ways to “remember” a subtree and copy it. One is of course by using variables; but their scope is restricted to a rule and therefore an arbitrary subtree can only be copied at a *fixed distance* of its original position. For instance in a rule of the form  $q(a(x)), D \rightarrow a(q_1(x), b(b(q_1(x))))$ , assuming that  $q_1$  copies its input, the copy of the original subtree is two-levels down from its next sibling but it cannot be arbitrary far. A second way to copy a subtree is to remember it using the states. Indeed, states can encode the knowledge that the TDTTR has accumulated along a path. However, since the number of states is finite, the only thing that a TDTTR can do is copy a fixed path. For instance for any given  $n$ , there exists a TDTTR that performs the transformation above, for trees of height  $n$  (it has essentially  $2^n - 1$  states which remember every possible path taken). For instance for  $n = 2$ , the TDTTR is:

$$\begin{aligned} q_0(a(x)), D &\rightarrow a(q_1(x)) \\ q_1(a()), D &\rightarrow a(a) \\ q_1(b()), D &\rightarrow a(b) \\ q_0(b(x)), D &\rightarrow b(q_2(x)) \\ q_2(a()), D &\rightarrow b(a) \\ q_2(b()), D &\rightarrow b(b) \end{aligned}$$

where  $\Sigma = \Delta = \{a, b\}$ ,  $Q = \{q_0, q_1, q_2\}$ ,  $Q_d = \{q_0\}$  and  $D = \{x \mapsto T_\Sigma\}$ . It is however impossible to write a TDTTR that replaces a leaf with a copy of the whole input, for inputs of arbitrary size. A similar example is used in [13] to show that TDTTR and bottom-up tree transducers are not comparable.  $\square$

## G. Operators on record types.

We use the theory of records defined for CDuce. We summarize here the main definitions. These are adapted from those given in Chapter 9 of Alain Frisch’s PhD thesis [15] where the interested



reader can also find detailed definitions of the semantic interpretation of record types and of the subtyping relation it induces, the modifications that must be done to the algorithm to decide it, finer details on pattern matching definition and compilation techniques for record types and expressions.

Let  $Z$  denote some set, a function  $r : \mathcal{L} \rightarrow Z$  is *quasi-constant* if there exists  $z \in Z$  such that the set  $\{\ell \in \mathcal{L} \mid r(\ell) \neq z\}$  is finite; in this case we denote this set by  $\text{dom}(r)$  and the element  $z$  by  $\text{def}(r)$ . We use  $\mathcal{L} \rightarrow Z$  to denote the set of quasi-constant functions from  $\mathcal{L}$  to  $Z$  and the notation  $\{\ell_1 = z_1, \dots, \ell_n = z_n, \_ = z\}$  to denote the quasi-constant function  $r : \mathcal{L} \rightarrow Z$  defined by  $r(\ell_i) = z_i$  for  $i = 1..n$  and  $r(\ell) = z$  for  $\ell \in \mathcal{L} \setminus \{\ell_1, \dots, \ell_n\}$ . Although this notation is not univocal (unless we require  $z_i \neq z$ ), this is largely sufficient for the purposes of this section.

Let  $\perp$  be a distinguished constant, then the sets  $\text{string} \rightarrow \text{Types} \cup \{\perp\}$  and  $\text{string} \rightarrow \text{Values} \cup \{\perp\}$  denote the set of all record types expressions and of all record values, respectively. The constant  $\perp$  represents the value of the fields of a record that are “undefined”. To ease the presentation we use the same notation both for a constant and the singleton type that contains it: so when  $\perp$  occurs in  $\mathcal{L} \rightarrow \text{Values} \cup \{\perp\}$  it denotes a value, while in  $\text{string} \rightarrow \text{Types} \cup \{\perp\}$  it denotes the singleton type that contains only the value  $\perp$ .

Given the definitions above, it is clear that the record types in Definition 2 are nothing but specific notations for some quasi-constant functions in  $\text{string} \rightarrow \text{Types} \cup \{\perp\}$ . More precisely, the open record type expression  $\{\ell_1:t_1, \dots, \ell_n:t_n, \dots\}$  denotes the quasi-constant function  $\{\ell_1 = t_1, \dots, \ell_n = t_n, \_ = \text{any}\}$  while the closed record type expression  $\{\ell_1:t_1, \dots, \ell_n:t_n\}$  denotes the quasi-constant function  $\{\ell_1 = t_1, \dots, \ell_n = t_n, \_ = \perp\}$ . Similarly, the optional field notation  $\{\dots, \ell?:t, \dots\}$  denotes the record type expressions in which  $\ell$  is mapped either to  $\perp$  or to the type  $t$ , that is,  $\{\dots, \ell = t \mid \perp, \dots\}$ .

Let  $t$  be a type and  $r_1, r_2$  two record type expressions, that is  $r_1, r_2 : \text{string} \rightarrow \text{Types} \cup \{\perp\}$ . The *merge* of  $r_1$ , and  $r_2$  with respect to  $t$ , noted  $\oplus_t$  and used infix, is the record type expression defined as follows:

$$(r_1 \oplus_t r_2)(\ell) \stackrel{\text{def}}{=} \begin{cases} r_1(\ell) & \text{if } r_1(\ell) \& t \leq \text{empty} \\ (r_1(\ell) \setminus t) \mid r_2(\ell) & \text{otherwise} \end{cases}$$

Recall that by Lemma 11 a *record type* (ie, a subtype of  $\{\dots\}$ ) is equivalent to a finite union of *record type expressions* (ie, quasi-constant functions in  $\text{string} \rightarrow \text{Types} \cup \{\perp\}$ ). So the definition of *merge* can be easily extended to all record types as follows

$$\left( \bigvee_{i \in I} r_i \right) \oplus_t \left( \bigvee_{j \in J} r'_j \right) \stackrel{\text{def}}{=} \bigvee_{i \in I, j \in J} (r_i \oplus_t r'_j)$$

Finally, all the operators we used for the typing of records in the rules of Section 4.2 are defined in terms of the merge operator:

$$t_1 + t_2 \stackrel{\text{def}}{=} t_2 \oplus_{\perp} t_1 \quad (1)$$

$$t \setminus \ell \stackrel{\text{def}}{=} \{\ell = \perp, \_ = c_0\} \oplus_{c_0} t \quad (2)$$

where  $c_0$  is any constant different from  $\perp$  (the semantics of the operator does not depend on the choice of  $c_0$  as long as it is different from  $\perp$ ).

Notice in particular that the result of the concatenation of two record type expressions  $r_1 + r_2$  may result for each field  $\ell$  in three different outcomes:

1. if  $r_2(\ell)$  does not contain  $\perp$  (ie, the field  $\ell$  is surely defined), then we take the corresponding field of  $r_2$ :  $(r_1 + r_2)(\ell) = r_2(\ell)$
2. if  $r_2(\ell)$  is undefined (ie,  $r_2(\ell) = \perp$ ), then we take the corresponding field of  $r_1$ :  $(r_1 + r_2)(\ell) = r_1(\ell)$

3. if  $r_2(\ell)$  may be undefined (ie,  $r_2(\ell) = t \mid \perp$  for some type  $t$ ), then we take the union of the two corresponding fields since it can results either in  $r_1(\ell)$  or  $r_2(\ell)$  according to whether the record typed by  $r_2$  is undefined in  $\ell$  or not:  $(r_1 + r_2)(\ell) = r_1(\ell) \mid (r_2(\ell) \setminus \perp)$ .

This explains all the examples we gave in the main text. In particular,  $\{a:\text{int}, b:\text{int}\} + \{a?:\text{bool}\} = \{a:\text{int} \mid \text{bool}, b:\text{int}\}$  since “ $a$ ” may be undefined in the right hand-side record while “ $b$ ” is undefined in it, and  $\{a:\text{int}\} + \{\dots\} = \{\dots\}$ , since “ $a$ ” in the right hand-side record is defined (with  $a \mapsto \text{any}$ ) and therefore has priority over the corresponding definition in the left hand-side record.

## H. Encoding of Co-grouping

As shown in Section 5.2, our `groupby` operator can encode JaQL’s `group each`  $x = e$  as  $y$  by  $e'$ , where  $e$  computes the grouping key, and for each distinct key,  $e'$  is evaluated in the environment where  $x$  is bound to the key and  $y$  is bound to the sequence of elements in the input having that key. Co-grouping is expressed by:

```
group
  l1 by x = e1 as y1
  ⋮
  ln by x = en as yn
into e
```

Co-grouping is encoded by the following composition of filters:

```
[ [l1] ... [ln] ];
[ Transform[x ⇒ (1, x)] ... Transform[x ⇒ (n, x)] ];
Expand;
groupby ((1, $) ⇒ [e1] | ... | (n, $) ⇒ [en]);
Transform[ (x, l) ⇒ ([l; Rgrp1] ... [l; Rgrpn]) ;
               [y1, ..., yn] ⇒ [e] ]
```

where

```
let filter Rgrpi = 'nil => 'nil
                  | ((i, x), tail) => (x , Rgrpi tail)
                  | _ => Rgrpi tail
```

Essentially, the co-grouping encoding takes as argument a sequence of sequences of values (the  $n$  sequences to co-group). Each of these  $n$  sequence, is tagged with an integer  $i$ . Then we flatten this sequences of tagged values. We can on this single sequence apply our `groupby` operator and modify each key selector  $e_i$  so that it is only applied to a value tagged with integer  $i$ . Once this is done, we obtain a sequence of pairs  $(k, l)$  where  $k$  is the common grouping key and  $l$  a sequence of tagged values. We only have to apply the auxiliary `Rgrp` filter which extracts the  $n$  subsequences from  $l$  (tagged with  $1..n$ ) and removes the integer tag for each value. Lastly we can call the recombining expression  $e$  which has in scope  $x$  bound to the current grouping key and  $y_1, \dots, y_n$  bound to each of the sequences of values from input  $l_1, \dots, l_n$  whose grouping key is  $x$ .

This figure "xhtml.png" is available in "png" format from:

<http://arxiv.org/ps/1303.1716v1>